

# Parametric Conditional Mean Inference with Functional Data applied to Lifetime Income Curves

JIN SEO CHO

Beijing Institute of Technology and Yonsei University

Email: jinseocho@yonsei.ac.kr

PETER C. B. PHILLIPS

Yale University, University of Auckland, University of Southampton & Singapore Management University

Email: peter.phillips@yale.edu

JUWON SEO

Department of Economics, National University of Singapore, AS2 #05-28, 1 Arts Link, Singapore 117570

Email: ecssj@nus.edu.sg

This version: August, 2020

## Abstract

We propose a framework for estimation of the conditional mean function in a parametric model with function space covariates. The approach employs a functional mean squared error objective criterion and allows for possible model misspecification. Under regularity conditions, consistency and asymptotic normality are established. The analysis extends to situations where the asymptotic properties are influenced by estimation errors arising from the presence of nuisance parameters. Wald, Lagrange multiplier, and quasi-likelihood ratio statistics are studied and asymptotic theory is provided. These procedures enable inference about curve shapes in the observed functional data. Several model specifications where our results are useful are analyzed, including time forms implied by panel data, random coefficient models, distributional mixtures, and copula mixture models. Simulations exploring the finite sample properties of our methods are provided. An empirical application conducts lifetime income path comparisons across different demographic groups according to years of work experience. Gender and education levels are found to produce differences in mean income paths, which reinforces earlier research results. But functional analysis reveals that the mean income paths are proportional so that, upon rescaling, the paths match over gender and across education levels.

**Key Words:** Functional data; Mean function; Wald test statistic; Lagrange multiplier test statistic; Quasi-likelihood ratio test statistic.

**Subject Class:** C11, C12, C80.

**Acknowledgements:** The acting Editor, Jesus Fernandez-Villaverde, and three anonymous referees provided very helpful comments on the original version of the paper for which we are most grateful. We also acknowledge helpful discussions with Kees Jan van Garderen, Kevin Sheppard, Richard Smith, Liangjun Su, Ying Wang, and participants of ANZESG (Wellington, 2019) and SETA (Osaka, 2019). Cho acknowledges research support from an Isaac Manasseh Meyer Fellowship of the National University of Singapore and kind hospitality of the Department of Economics at the Chinese University of Hong Kong during his visit in 2020; Phillips acknowledges research support from a Kelly Fellowship at the University of Auckland and the NSF under Grant No. SES 18-50860; and Seo acknowledges research support from AcRF Tier 1.

# 1 Introduction

Functional data analysis (FDA) has been attracting increasing attention in the statistical and econometric literatures. Among the reasons for this growing interest are the growing availability of very large cross section and spatio-temporal datasets, the inherent interest in studying function space, curve, or surface realizations of data, and the potential that function space methods provide for economic, financial, and scientific data analysis at this enhanced level of detail. In place of individual point observations, functional data lead naturally to the analysis of continuous phenomenon such as time series curves that record trend or growth trajectories and do so under assumptions that can allow great generality. Ramsay and Dalzell (1991), Rice and Silverman (1991), Ramsay and Silverman (1997), Bosq (2000), and Horvath and Kokoska (2012) are now classical references on FDA. We also refer to Ramsay and Silverman (2002), Cai and Hall (2006), Ferraty and Vieu (2006), Cardot, Crambes, Kneip, and Sarda (2007), Hall and Horowitz (2007), Zhang and Chen (2007), Müller, Sen, and Stadtmüller (2011), Cao, Yang, and Todem (2012), and Müller (2012) for further background and recent research.

Theoretical developments in FDA often focus on nonparametric model analyses such as functional regression models, so that the objects of interest in estimation and inference are typically nonparametric functions or operators acting on a function space in which the data are defined. Such analyses provide a useful foundation for a general approach to FDA. But the final objects in fully nonparametric studies may sometimes be difficult to interpret in practice. Rather than pursue a nonparametric approach, the present study works with a parametric formulation that aims to be amenable to implementation and interpretation in application.

For this goal, we propose a novel and efficient framework for the estimation and inference of the (conditional) mean function of functional data. Our approach is different from many other recent FDA studies in that we assume a parametric model for the mean function, one that may possibly be misspecified, whereas the observations remain nonparametric random elements in a measurable function space. We study the influence of potential misspecification of the mean function on our estimator by examining it in parallel to the analysis of a quasi-maximum likelihood estimator, much as in White's (1982, 1994) investigations with finite dimensional data.

Our approach has some advantages that are useful in practical research. First, it allows us to construct simple statistical tests for the (conditional) mean functions with the estimated parameters using the asymptotic machinery. In the nonparametric context, inference about the slope function in functional regression often involves technically challenging issues arising from well-known ill-posed functional inversion problems. In our framework on the other hand, the relevant null hypothesis can be easily tested by estimating vector valued parameters, enabling straightforward inference on the mean function. Our approach also provides a convenient way to study the derivatives of population mean functions. Derivative functions such as growth rates of income, wealth, or employment are often of significant interest in economic applications. With nonparametric approaches numerical differentiation methods are typically employed, whereas parametric model estimation enables estimation and testing for exact analytic derivatives. For example, statistical analysis of functional data such as shift registration alignment requires the estimation of the mean function and its exact derivatives, for which we can conveniently apply the methods developed herein.

Some well-known statistical methods that are now used in econometrics provide further motivation for the present parametric functional data framework. For instance, there is a large literature in meta analysis on the combination of independent test statistics, some methods using weighting techniques, in which  $p$ -curve analysis may be used. These methods originated in the work of Fisher (1932) (see also Pearson, 1950; Lancaster, 1961; van Zwet and Oosterhoff, 1967; Westberg, 1985, among others) and are now employed - assess selective reporting in empirical work in terms of publication bias and  $p$ -hacking - e.g., see Simonsohn, Nelson and Simmons (2014). Such methods can be viewed in terms of functional data analysis with a

particular form of weight function. Other relevant literatures are the unidentified model analysis of Davies (1977, 1987) and the minimum distance tests in Pollard (1980). Both provide motivation for the present approach in which the tests are interpreted as statistics obtained from functional observations – see Section 5 of this study for an example. There is now a substantial literature extending the inferential methods of Davies (1977, 1987) and Pollard (1980) and our approach can be applied to similar problems with various identification features or empirical distributions (e.g., Hansen, 1996; Andrews, 2001; Baek, Cho, and Phillips 2015; Bierens, 1990; Cho, Park, and Phillips 2018; Cho and White, 2007, 2010, 2018b; Stinchcombe and White, 1998, and the references therein). Each observation that underlies the statistics in these studies is a functional observation, making the present approach applicable.

The literature on FDA is growing and the direction of the present research contributes to that expanding literature. To mention a few recent developments, we draw attention to the following studies. Grenander (1981) and Kutoyants (1984) study the estimation of a parametric model for a single functional observation. Crambes, Gannoun, and Henchiri (2013) examine the estimation of a quantile regression function with a functional covariate by means of a support vector machine (SVM) learning method. Here the dependent variable is a real random variable and the explanatory variable is a functional observation. The authors first apply a linear integral operator to the functional data, converting the functional observation to a random variable, thereby enabling estimation of the conditional quantile functional (i.e., the quantile function between the dependent variable and the transformed observations) using SVM methodology. Zhang and Chen (2007) examine the so-called “smoothing first, then estimation” principle which substitutes functional observations obtained by applying local polynomial kernel estimation to discrete data to estimate the unobserved underlying functional observations. Under some mild regularity conditions, they show that the influence of this substitution on inference can be ignored as the sample size increases. Li, Robinson, and Shang (2020) study time series of function space curves under long range dependence, establishing limit theory for sample averages, estimating the covariance kernel function of the functional data via functional principal component analysis, and using orthonormal functions to span the dominant subspace of the curves. Chang, Hu, and Park (2019) consider estimation of a functional autoregressive model with serially correlated functional data and establish consistency and asymptotic normality of the autoregressive operator estimator. Phillips and Jiang (2019) study parametric autoregression with function valued time series in stationary and nonstationary cases, establish asymptotic theory of estimation and inference, and apply the methods to analyze household Engel curves among ageing seniors in a wide panel dataset. These papers all relate to the current study in terms of the use of functional data but differ from its focus on estimation and inference of a parametric nonlinear conditional mean function involving functional observations and given vector valued explanatory variables, allowing for possible misspecification.

The paper is organized as follows. Section 2 describes the data, defines the model, and sets up the estimation criterion in terms of a functional mean squared error (FMSE). Section 3 proposes a functional least squares (FLS) estimator for the parametric conditional mean function. Consistency and asymptotic normality of the FLS estimator is established allowing for possible nuisance effects. Asymptotic covariance matrix estimation is also discussed. Section 4 provides a general framework for hypothesis testing, with extensions of Wald, Lagrange multiplier (LM), and quasi-likelihood ratio (QLR) test statistics and asymptotics to the functional data environment. In Section 5, several model specifications where our results are useful are analyzed, including panel data inference, random coefficient models, distributional mixtures, and copula mixture models. Finite sample simulations are also provided. Section 6 reports an empirical application of the methods to lifetime income path comparisons across different demographic groups. Conclusions are given in Section 7 and proofs are collected in the Appendix.

## 2 Setup

We are interested in studying data that comprise a set of observable random functions  $g_i(\cdot)$  and observable random vectors  $x_i$ , which are given as

$$\{(g_i(\cdot), x_i) : g_i : \Gamma \mapsto \mathbb{R} \text{ and } x_i \in \mathbb{R}^k\}_{i=1}^n, \quad (1)$$

where  $n$  is the sample size and  $k \in \mathbb{N}_+$ .

There are many examples fitting this data structure. We provide below a list of some of the data sets with this structure to help fix the main ideas of the present paper:

- In our empirical work the function  $g_i$  is used to represent an observable curve that shows an individual's income profile over their lifetime working years or some relevant subset of those years, such as those that follow 10 or more years work experience, signifying maturity in the labor force. The vector  $x_i$  embodies relevant individual characteristics such as gender and educational level attained which influence the earnings function  $g_i$ . The curve  $g_i(\gamma)$  with  $\gamma \in \Gamma = [10, 40]$  then represents the earnings profile of individual  $i$  with those particular characteristics over a lifetime of 30 years with at least 10 years work experience. The primary object of interest is then the conditional mean function of the income profile curve given the observable characteristics  $x$ . This function might then be modeled using a specific parametric functional form in terms of work experience for those given characteristics, a quadratic function being  $E_{\mathbb{P}}[g_i(\gamma)|x_i = x] = m_1(\theta, x) + m_2(\theta, x)\gamma + m_3(\theta, x)\gamma^2 =: \mu(\gamma; \theta, x)$  for pre-specified functions  $m_j(\theta, x)$  ( $j = 1, 2, 3$ ), where  $\theta$  is a parameter vector to be estimated from the functional data  $g_i(\cdot)$  and the relevant characteristics  $x_i$ . The notation  $E_{\mathbb{P}}[g_i(\gamma)|x_i = x]$  indicates that expectation is taken with respect to the probability measure  $\mathbb{P}$ , defined in Assumption 1, conditional on the given characteristics  $x$ , as explained in (2) below. There are many mean functions in the literature similar to the income profile. For example, the  $S$ -shaped curves or sigmoid functions are popularly employed as mean functions of business growth, crop yield, and learning outcome data, among others, and sigmoid functions such as the logistic, arctangent, and error functions are often used in empirical work. As another example, trigonometric functions are commonly used in the literature to capture various cyclical patterns.
- An observed measurable functional transformation of a random element can be treated as a functional observation. Specifically, we may consider a random variable  $(y_i, x_i)$  and its parametric transformation via a measurable function:  $g(\gamma, y_i, x_i)$ , where  $\gamma \in \Gamma \subset \mathbb{R}^d$  as in our paper. For simplicity, we let  $g_i(\cdot) := g(\cdot, y_i, x_i)$  be a functional observation defined on  $\Gamma$  and analyze its (conditional) mean function  $\mathbb{E}[g_i(\cdot)|x_i]$ , defined according to the researcher's interest. For example, an empirical researcher may wish to test a particular form of the conditional mean function that may be implied by an hypothesis of interest. The examples in Sections 5 handle cases in which the mean function is identically zero (e.g., from Davies' (1977, 1987) identification problem) and the hypotheses of interest are imposed on the score function, which is now the parametric transformation. For testing purposes, the paper develops methodologies to test the hypothesized mean function form by extending traditional Wald, Lagrange multiplier and likelihood ratio test principles to functional data.
- The analysis of  $p$ -curves can be interpreted in terms of functional data. If we let  $g_i(\gamma) := K_h(\gamma - x_i)$ , where  $K_h(\cdot)$  is a kernel function with bandwidth  $h$ , the sample average of the  $g_i(\cdot)$  becomes a standard kernel density function using functional observations. The mean function of  $g_i(\cdot)$  conveys informative characteristics of the data and  $p$ -curve analysis examines the density function shape of the  $p$ -value. Under the null hypothesis the  $p$ -value of a test statistic follows a uniform distribution. So if the mean function of  $g_i(\cdot)$  is estimated by the kernel density function over the region around a test level such as 0.05, the shape should be flat under the null or a decreasing function under the alternative, which is a

prime object of investigation under  $p$ -curve analysis.

- Functional data are often obtained by interpolation methods. Individual observations obtained at discrete locations and/or times typically display different random patterns and unobservable observations between discrete locations and/or times may be predicted by local polynomial kernel, sieve, kriging, and polynomial spline estimation, among other methods, thereby converting discrete observations into functional data (e.g., Zhang and Chen, 2007; Chen, 2007; Chilès and Delfiner, 1999; Wang, 2011). The (conditional) mean function estimated from the resulting functional data can be exploited to infer properties of the (conditional) mean function of the missing data between the discrete observations. For panel data sets with discretely collected observations functional data delivered by such interpolations can be exploited to infer (conditional) mean equations, although inference may not be straightforward if unbalanced panel data are used.

Many more functional data examples fit the modeling structure employed in the present study. We focus on the first two examples above to demonstrate the proposed methodology and employ these as our running examples.

The following conditions provide a formal probabilistic framework in which data of this type may be analyzed.

**Assumption 1. (Data):** (i)  $(\Omega, \mathcal{F}, \mathbb{P})$  is a complete probability space and  $(\Gamma, \tau)$  is a compact metric space with  $\tau$  being a metric defined on  $\Gamma$ ;

(ii)  $\{(g_i(\cdot, \cdot), x_i')' : g_i : \Omega \times \Gamma \mapsto \mathbb{R} \text{ and } x_i : \Omega \mapsto \mathbb{R}^k\}_{i=1}^n$  is a set of identically and independently distributed (IID) observations such that for each  $\gamma \in \Gamma$ ,  $\{(g_i(\cdot, \gamma), x_i')'\}$  is  $\mathcal{F}$ -measurable, and  $g_i(\omega, \cdot) \in \mathcal{C}^{(0)}(\Gamma)$  for all  $\omega \in F \in \mathcal{F}$  with  $\mathbb{P}(F) = 1$ , where  $\mathcal{C}^{(\ell)}(\cdot)$  denotes the space of  $\ell$ -times continuously differentiable functions;

(iii)  $(\Gamma, \mathcal{G}, \mathbb{Q})$  and  $(\Omega \times \Gamma, \mathcal{F} \otimes \mathcal{G}, \mathbb{P} \times \mathbb{Q})$  are complete probability spaces, and  $g_i(\cdot, \cdot)$  is  $\mathcal{F} \otimes \mathcal{G}$ -measurable.  $\square$

The argument space  $\Gamma$  is the space where the functional observations are defined for a fixed  $\omega \in \Omega$ . For convenience, we define the functional observations  $g_i$  on the product space of  $\Omega \times \Gamma$  rather than interpreting them as elements in some Hilbert space. A straightforward example is when  $g(\omega, \cdot)$  is a function defined on a subset of  $s$ -dimensional Euclidean space ( $s \in \mathbb{N}_+$ ). For such a case,  $\Gamma$  and  $\tau$  are the subset and the Euclidean distance, respectively. In Assumption 1(iii), the original probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  is combined with  $(\Gamma, \mathcal{G}, \mathbb{Q})$  to form the product probability space  $(\Omega \times \Gamma, \mathcal{F} \times \mathcal{G}, \mathbb{P} \times \mathbb{Q})$ . We call  $(\Gamma, \mathcal{G}, \mathbb{Q})$  and  $\mathbb{Q}$  the *adjunct probability space* and *adjunct probability measure*, respectively. We suppose that these spaces are judiciously chosen by the researcher to match the specific modeling interests and goals.

To avoid confusion this formulation does *not* mean that  $\gamma$  is an unobserved variable with a known distribution, for which we may apply a minimum distance estimation technique to infer the data structure. Our model and approach are therefore fundamentally different from earlier general work on conditional moment estimation, such as the work of Ai and Chen (2003). On the contrary, the adjunct probability space and adjunct probability measure are devices that are carefully selected by the researcher for the empirical purpose in hand and these may therefore differ according to the study. This important feature of our approach means that the analysis of functional data is conducted quite differently from formulations in which  $\gamma$  is assumed to be a random variable, such as in the context of models with conditional moment restrictions containing observed functions but unobserved random variables. Furthermore, the carefully selected devices may in turn be influenced directly by properties associated with estimators of the conditional mean function and test statistics used for inference. As we detail at the end of Section 3.2, the powers of test statistics devised for inference on the conditional mean function can be affected by the selection of the adjunct probability measure. This fact may be used by investigators to guide the choice of  $\mathbb{Q}$  judiciously so that tests employed may be conveniently computed and have power against what are viewed as more realistic alternatives.

The adjunct probability space and adjunct probability measure are not directly part of the stochastic aspects of the data, as

mentioned above. But the measurability conditions employed in Assumption 1 are useful in defining some of the integrals that are introduced below. If  $g_i(\cdot)$  is continuous in  $\gamma$  almost surely (a.s.) with respect to  $\mathbb{P}$ , the joint measurability condition in Assumption 1(iii) trivially holds by lemma 2.15 of Stinchcombe and White (1992).<sup>1</sup> Otherwise, the measurability condition has to be verified by explicitly considering the properties of the function. It will be convenient to proceed in this development by requiring the continuity condition to hold, along with other conditions for the investigator to consider when choosing the adjunct probability measure for a regularly behaving FLS estimator. Also, in the notation (1) and elsewhere, the argument  $\omega$  is often suppressed for convenience.

We further suppose that the primary subject of interest is the conditional mean function of  $g_i$ , which is defined by the integral<sup>2</sup>

$$\mu(\gamma; x) := \int g(\gamma) d\mathbb{P}(g(\gamma)|x), \quad (2)$$

where  $\mathbb{P}(\cdot|x)$  is the conditional probability measure of  $g_i(\gamma)$  given  $x_i = x$ . For each  $\gamma \in \Gamma$ , we treat the function  $g(\gamma)$  as a random variable and compute its conditional mean  $\mu$ . Therefore, if we let  $E_{\mathbb{P}}$  denote the expectation operator associated with the probability measure  $\mathbb{P}$ ,  $\mu(\gamma, x)$  can be expressed as  $E_{\mathbb{P}}[g_i(\gamma)|x_i = x]$ . If the function  $g_i(\cdot)$  is constant a.s., we can view it as a simple random variable, so that  $E_{\mathbb{P}}[g_i(\gamma)|x_i = x]$  becomes the conventional conditional mean of  $g_i(\cdot)$ .

For a parametric specification of the conditional mean with parameter vector  $\theta$  we may write  $\mu(\gamma; x) = \mu(\gamma; \theta, x)$ , as earlier. That is,  $\mu(\gamma; x)$  is a special case of  $\mu(\gamma; \cdot, x)$  that can be obtained by plugging  $\theta$  into the empty argument space. More generally, we wish to allow for misspecification in our analysis. To do so, we define  $\mathcal{M}$  to be a collection of parametric models specified by a function  $\rho$ . Specifically, for each  $x$ ,

$$\mathcal{M} := \{\rho(\cdot, \theta, x) : \Gamma \mapsto \mathbb{R} | \theta \in \Theta \subset \mathbb{R}^d\}.$$

The following conditions are assumed for  $\mathcal{M}$  and  $\rho$ :<sup>3</sup>

**Assumption 2. (Model):** (i) For each  $\theta \in \Theta$ ,  $\rho(\cdot, \theta, \cdot) : \Gamma \times \Omega \mapsto \mathbb{R}$  is  $\mathcal{F} \otimes \mathcal{G}$ -measurable, where the parameter space  $\Theta$  is a compact and convex set in  $\mathbb{R}^d$  for  $d \in \mathbb{N}_+$ ;

(ii) for each  $\gamma \in \Gamma$ ,  $\rho(\gamma, \cdot, \omega) : \Theta \mapsto \mathbb{R} \in \mathcal{C}^{(2)}(\Theta)$  for all  $\omega \in F \in \mathcal{F}$  with  $\mathbb{P}(F) = 1$ ;

(iii) for each  $\theta \in \Theta$ ,  $\rho(\cdot, \theta, \omega) \in \mathcal{C}^{(0)}(\Gamma)$  for all  $\omega \in F \in \mathcal{F}$  with  $\mathbb{P}(F) = 1$ ; and

(iv) the optimizer  $\theta_*$  is unique and lies in the interior of the parameter space  $\Theta$ , where  $\theta_* := \arg \min_{\theta \in \Theta} q(\theta)$  and  $q(\theta) := \int \int \{g(\gamma) - \rho(\gamma, \theta, x)\}^2 d\mathbb{P}(g(\gamma), x) d\mathbb{Q}(\gamma)$ . □

Based on the specification of the mean function  $\mu$ , the functional mean squared error (FMSE) criterion for estimation of the parametric mean function is defined by the functional  $q(\cdot)$  in analogy to the usual mean squared error (MSE) in least squares estimation. As in standard analysis (iv) requires a unique optimizer  $\theta_*$  in the interior of  $\Theta$ , thereby avoiding possible non-identified model issues from consideration in this development. The identification condition is imposed as a bottomline condition to assist the primary goal of this study to develop functional data analysis in parallel to standard regression analysis with possibly

<sup>1</sup>Even if  $g_i(\cdot)$  is discontinuous, as it is for the score function derived from testing for a structural break without knowing the break point, an analysis similar to that of the current study can be employed with some modifications to allow for broken functional data, either by using a smudge function in place of the indicator break which can be approximated by a local polynomial kernel, or by extending the methodology directly to allow directly for indicator function breaks. Analysis of this topic is left for future research.

<sup>2</sup>Here and throughout the rest of the paper, unless otherwise noted, we use notations such as (2) to denote integrals over the whole probability space, so that here  $\int g(\gamma) d\mathbb{P}(g(\gamma)|x) = \int_{\Omega} g(\omega, \gamma) d\mathbb{P}(g(\gamma)|x)$ .

<sup>3</sup>Note that this framework is significantly different from Bugni, Hall, Horowitz, and Neumann (2009) which is concerned with parametric specifications of functional observations. In particular, our formulation means that we can let  $g_i$  be a random function with no further specification, although its mean function is parametrically specified.

misspecified models (e.g., White, 1982, 1994; Newey and McFadden, 1994). The model identification condition plays a central role in regression analysis with misspecified models when proving convergence and deriving limit theory, letting the estimated model be a close approximate to the conditional mean, even if the model is misspecified. As our model exercises and empirical analysis demonstrate in Sections 5 and 6, the empirical researcher has freedom to specify  $\rho(\cdot, x)$  and  $\mathbb{Q}$  to ensure they satisfy Assumption 2(iv) and do not suffer from unidentified model issues.

Several additional technical conditions on  $g_i$  and  $\rho(\cdot, \cdot, x_i)$  are given in the following assumption to assist in subsequent derivations. In what follows, for a probability measure  $\mathbb{P}$  on  $\Omega$ , let  $L^\ell(\mathbb{P}) := \{f : \int_\Omega |f(\omega)|^\ell d\mathbb{P}(\omega) < \infty\}$ , for  $\ell = 1$  and  $2$ .

**Assumption 3. (Moments):** For some  $m_i \in L^2(\mathbb{P})$ ,

(i)  $\sup_{\gamma \in \Gamma} |g_i(\gamma)| \leq m_i$  a.s.  $-\mathbb{P}$ ;

(ii)  $\sup_{(\gamma, \theta) \in \Gamma \times \Theta} |\rho(\gamma, \theta, x_i)| \leq m_i$  a.s.  $-\mathbb{P}$ ;

(iii) for each  $j = 1, 2, \dots, d$ ,  $\sup_{(\gamma, \theta) \in \Gamma \times \Theta} |(\partial/\partial\theta_j)\rho(\gamma, \theta, x_i)| \leq m_i$  a.s.  $-\mathbb{P}$ ;

(iv) for each  $j$  and  $j' = 1, 2, \dots, d$ ,  $\sup_{\theta \in \Theta} |(\partial^2/\partial\theta_j\partial\theta_{j'})\rho(\cdot, \theta, x_i)| \leq m_i$  a.s.  $-\mathbb{P}$ . □

Assumption 3 is imposed to ensure by domination the existence of  $q(\cdot)$  and a global minimum of this functional in conjunction with Assumption 2(iv). The moment conditions (iii) and (iv) also ensure that first and second order conditions apply for minimization of  $q(\cdot)$ .

In practice, the functional form of  $\mu$  is typically unknown and the given model class  $\mathcal{M}$  may not contain a parameter value  $\theta$  such that  $\rho(\cdot, \theta, x) = \mu(\cdot, x)$  for all  $x$ . We say  $\mathcal{M}$  is *correctly specified* if there exists a parameter value  $\theta_0 \in \Theta$  such that  $\mu(\cdot, x_i) = \rho(\cdot, \theta_0, x_i)$  a.s.  $-\mathbb{P} \cdot \mathbb{Q}$ . Otherwise, we say that  $\mathcal{M}$  is *misspecified*, in which case  $\theta_0$  is undefined. Theorem 1 below provides a useful decomposition of the functional  $q(\theta)$  that characterizes the implications of correct specification and misspecification on the minimizer of  $q(\cdot)$ .

**Theorem 1.** Given Assumptions 1, 2, and 3, we have

$$q(\theta) = \int \int \text{var}_{\mathbb{P}}[g_i(\gamma)|x] d\mathbb{P}(x) d\mathbb{Q}(\gamma) + \int \int \{\mu(\gamma, x) - \rho(\gamma, \theta, x)\}^2 d\mathbb{P}(x) d\mathbb{Q}(\gamma),$$

where for each  $\gamma$ ,  $\text{var}_{\mathbb{P}}[g_i(\gamma)|x] := \int \{g(\gamma) - \mu(\gamma, x)\}^2 d\mathbb{P}(g(\gamma)|x)$ . □

When  $\mathcal{M}$  is correctly specified, we have  $q(\theta_0) = \int \int \text{var}_{\mathbb{P}}[g_i(\gamma)|x] d\mathbb{P}(x) d\mathbb{Q}(\gamma)$  with  $\theta_* = \theta_0$ , so that the mean function  $\mu$  is uniquely identified. In this case, the FMSE cannot be smaller than  $\int \int \text{var}_{\mathbb{P}}[g_i(\gamma)|x] d\mathbb{P}(x) d\mathbb{Q}(\gamma)$ . On the other hand, when  $\mathcal{M}$  is not correctly specified,  $\theta_0$  cannot be identified by minimizing  $q(\cdot)$ . This is because the FMSE is affected by an additional term that reflects the error impact of model misspecification. In this case,  $\theta_*$  should be understood as a parameter value of  $\theta$  which minimizes the sum of two squared errors, one being the mean squared error obtained if the model had been correctly specified, the other being the squared error component arising from model misspecification. In general, we may not presume that the model  $\mathcal{M}$  is correctly specified unless additional information on the parametric form of the mean function is provided. This presumption stems from the fact that empirical models are easily misspecified. In spite of considerable effort and diagnostic testing empirical models are typically only approximations and the same is true for models with functional data. Our treatment allows for  $\mathcal{M}$  to be a possibly misspecified model class and asymptotic properties are developed in this setting. The framework for the asymptotic theory then provides a wider understanding of the estimation procedures, thereby enabling a study of the impact of misspecification and opening up the possibility to detect model misspecification using devices such as the information matrix test. This research direction is only mentioned and is not pursued in detail in the present paper.

In what follows, it is convenient to employ a slight abuse of notation and use  $\mu_i(\cdot)$  and  $\rho_i(\gamma, \theta)$  to denote  $\mu(\cdot, x_i)$  and  $\rho(\gamma, \theta, x_i)$ , respectively. We also abbreviate  $\int g(x)dF(x)$  and  $\int \int k(x, y)dF(x, y)$  as  $\int g(x)dF$  and  $\int \int k(x, y)dF$ , respectively.

### Example 1: Distributional Specification

We introduce a mixture model illustration that is continued as a running example in the paper to illustrate our methods. A second example is provided below by accommodating a nuisance parameter estimation error effect in the FDA. These examples are employed in the simulations in Section 5 along with other models.

Finite mixture models are popular for constructing more flexible distribution functions and for modeling clustered data. They can also be used for certain types of distribution specification tests. To fix ideas, let  $f_i(\cdot; \theta_{i\uparrow})$  be a component density function for  $i = 1, \dots, K$ . The component densities can be chosen from the same family or from different families of distributions. Accordingly, the parameter vector  $\theta_{\uparrow} = (\theta_{1\uparrow}, \theta_{2\uparrow}, \dots, \theta_{K\uparrow})'$  is defined on the product of each parameter space,  $\Theta_1 \times \Theta_2 \times \dots \times \Theta_K$ . For weights  $\pi_i \in [0, 1]$  with  $\sum_{i=1}^K \pi_i = 1$ , the corresponding finite mixture model (Everitt and Hand, 1981; McLachlan and Peel, 2004; Schlattmann, 2009) is defined as the weighted sum

$$f(\cdot; \pi_1, \dots, \pi_K, \theta) = \sum_{i=1}^K \pi_i f_i(\cdot; \theta_i).$$

Various types of parametric distributions have been explored in making such constructions in the literature. For example, mixtures of normals, binomials, gammas and von Mises distributions are popular in applications. Amongst many references in the literature to the use of mixtures we mention Chernoff and Lander (1995), Liang and Rathouz (1999), Chen, Chen, and Kalbfleisch (2001), Cho and White (2007, 2010), Fu, Chen, and Yi (2008), Chen and Li (2009), Ning, Gupta, Yu, and Zhang (2009), Niu, Li, and Zhang (2011), and Wong and Li (2014).

We focus here on inference concerning sample homogeneity as an application of distribution specification tests. For a specific example consider the following mixture of exponential distributions

$$f(x; \pi_{\uparrow}, \gamma_{\uparrow}) = (1 - \pi_{\uparrow}) \exp(-x) + \pi_{\uparrow} \gamma_{\uparrow} \exp(-\gamma_{\uparrow} x),$$

where  $\gamma_{\uparrow} \in \Gamma := [\underline{\gamma}, \bar{\gamma}]$ . For simplicity, assume  $\underline{\gamma} > 1$  and  $\bar{\gamma} < \infty$ .<sup>4</sup> A primary concern in this mixture distribution is to test whether  $\pi_{\uparrow} = 0$  for if this hypothesis holds the observations are believed to come from a homogeneous population that follows a standard exponential distribution. Davies (1977) applied Neyman's (1959)  $C(\alpha)$  test principle and derived a maximal test statistic defined as

$$\sup_{\gamma \in \Gamma} n^{-1/2} \sum_{i=1}^n g_i(\gamma), \quad \text{where} \quad g_i(\gamma) := \frac{(2\gamma - 1)^{1/2}}{\gamma - 1} \{\gamma \exp[(1 - \gamma)x_i] - 1\},$$

and  $\{x_i\}$  is a set of IID observations. Note that in this formulation  $g_i(\cdot)$  is a random function defined on  $\Gamma$  obtained by standardizing the score with respect to  $\pi$  and evaluating it at  $\pi = 0$ , so that we may treat it as a functional observation and apply the theory of the present study.

If the assumed mixture model is correct, the population mean function of  $g_i(\cdot)$  can be accordingly derived. Note that under

<sup>4</sup>The definition of  $g_i(\cdot)$  implies that  $g_i(\cdot)$  is not defined for  $\gamma < 1/2$ , and  $g_i(1)$  is not defined, although  $\lim_{\gamma \rightarrow 1} g_i(\gamma)$  can be obtained by L'hôpital's rule, giving  $\lim_{\gamma \rightarrow 1} g_i(\gamma) = -(x_i - 1)$ . Our model exercise avoids this feature by setting  $\underline{\gamma} > 1$ .



the DGP described above, for each  $\gamma$ ,  $\mu(\gamma)$  is obtained as

$$\mathbb{E}[g_i(\gamma)] := \mu(\gamma) = \pi_{\dagger} \frac{(\gamma_{\dagger} - 1)(2\gamma - 1)^{1/2}}{(\gamma + \gamma_{\dagger} - 1)},$$

implying that  $\mu(\cdot) \equiv 0$  if  $\pi_{\dagger} = 0$ , whereas  $\mu(\cdot)$  is a non-zero function of  $\gamma$  with unknown parameter  $\gamma_{\dagger}$ , motivating  $C(\alpha)$  test statistic by finding  $\gamma$  such that  $\mu(\gamma) > 0$ , from which the power of the test is acquired. That is,  $\gamma_{\dagger}$  is identified only when  $\pi_{\dagger} \neq 0$ . Nonetheless, the null limit distribution of  $C(\alpha)$  test depends on the covariance kernel of a Gaussian process obtained in Section 5, and this makes its application inconvenient. Furthermore, if the mixture assumption is wrong, the application of  $C(\alpha)$  test is not statistically precise as its null limit distribution is obtained by the misspecified model.

Motivated by this feature, we desire to provide straightforwardly applicable diagnostic statistics for testing the functional form of  $\mu(\cdot)$ . For this goal we first specify a model for  $\mu(\cdot)$ . The functional form of  $\mu(\cdot)$  indicates a particular direction for model specification. If  $\gamma$  is close to  $\gamma_{\dagger}$ ,  $\mu(\gamma)$  is approximable by  $(\gamma - 1)/(2\gamma - 1)^{1/2}$ , enabling us to specify a simple model for  $\mu(\cdot)$  by

$$\rho(\gamma, \theta_1, \theta_2) := \theta_1 \rho_1(\gamma) + \theta_2 \rho_2(\gamma) := \theta_1 + \theta_2 \frac{(\gamma - 1)}{(2\gamma - 1)^{1/2}}$$

in the sense that it is linear with respect to  $(\theta_1, \theta_2)$  with two functional regressors  $\rho_1(\cdot) := 1$  and  $\rho_2(\cdot) := ((\cdot) - 1)/(2(\cdot) - 1)^{1/2}$ . Note that the model  $\mathcal{M}$  equipped with  $\rho(\cdot)$  is misspecified for  $\mu(\cdot)$  even if the mixture assumption is correct for  $x_i$ , but  $(\theta_{1*}, \theta_{2*})$  has to equal  $(0, 0)$  for  $\pi_{\dagger} \neq 0$  as shown in the following paragraph. In contrast,  $(\theta_{1*}, \theta_{2*}) = (0, 0)$ , if  $\mu(\cdot) \equiv 0$  for whatever reason. This simple fact motivates estimation of  $(\theta_{1*}, \theta_{2*})$  and testing whether it is zero.

Before moving to the estimation of  $(\theta_{1*}, \theta_{2*})$ , we verify the regularity conditions. First, the conditions in Assumption 1 are trivially satisfied from the condition that  $\Gamma$  is a compact set in the Euclidean real line; each observation is IID; and  $g_i(\cdot)$  is a continuous function with probability 1. Next,  $\rho(\cdot)$  is continuous on  $\Gamma$  for each  $(\theta_1, \theta_2)$  and linear with respect to  $(\theta_1, \theta_2)$  for each  $\gamma$ , satisfying Assumption 2 by further supposing that  $(\theta_{1*}, \theta_{2*})$  is an interior element of a compact and convex parameter space  $\Theta$  in  $\mathbb{R}^2$ , where  $(\theta_{1*}, \theta_{2*})$  is obtained by minimizing  $q(\cdot)$  in Theorem 1:

$$\begin{bmatrix} \theta_{1*} \\ \theta_{2*} \end{bmatrix} = \begin{bmatrix} \int \rho_1^2(\gamma) d\mathbb{Q}(\gamma) & \int \rho_1(\gamma) \rho_2(\gamma) d\mathbb{Q}(\gamma) \\ \int \rho_1(\gamma) \rho_2(\gamma) d\mathbb{Q}(\gamma) & \int \rho_2^2(\gamma) d\mathbb{Q}(\gamma) \end{bmatrix}^{-1} \begin{bmatrix} \int \rho_1(\gamma) \mu(\gamma) d\mathbb{Q}(\gamma) \\ \int \rho_2(\gamma) \mu(\gamma) d\mathbb{Q}(\gamma) \end{bmatrix}.$$

Here, the integrals are computed over  $\Gamma$ , and the inverse matrix exists by Cauchy-Schwarz, implying that  $(\theta_{1*}, \theta_{2*})$  is well defined and unique. Here, if  $\mu(\cdot) \equiv 0$ , it follows  $(\theta_{1*}, \theta_{2*}) = (0, 0)$ . If the mixture assumption is further assumed to be correct without assuming that  $\pi_{\dagger} = 0$ ,  $(\theta_{1*}, \theta_{2*})$  is computed by imposing the functional form of  $\mu(\cdot)$  implied by the mixture assumption. That is, if we let  $\Gamma := [1.5, 2.5]$  and  $\gamma_{\dagger} = 2$  as for the simulations in Section 5, it follows that  $(\theta_{1*}, \theta_{2*}) = \pi_{\dagger}(0.5687, 0.0101)$  by applying Gauss-Legendre's numerical quadrature, so that if  $\pi_{\dagger} = 0$ , then  $(\theta_{1*}, \theta_{2*}) = (0, 0)$ . Finally, note that  $|g_i(\cdot)| < \bar{\gamma}(2\underline{\gamma} - 1)^{1/2}/(\underline{\gamma} - 1)$ , which implies that it is trivial to find  $m_i$  satisfying Assumption 3(i). Furthermore, Assumptions 3(ii-iii) also trivially hold by the linearity of  $\rho(\cdot)$  with respect to  $(\theta_1, \theta_2)$  and the compact space condition for  $\Theta$  and  $\Gamma$ .  $\square$

### 3 Functional Least Squares (FLS)

#### 3.1 FLS Estimation without Nuisance Effects

We consider estimation of  $\theta$  based on the FMSE criterion. The *functional least squares* (FLS) parametric estimator is defined as the extremum estimator

$$\hat{\theta}_n := \arg \min_{\theta \in \Theta} q_n(\theta), \quad \text{where} \quad q_n(\theta) := \frac{1}{n} \sum_{i=1}^n \int \{g_i(\gamma) - \rho_i(\gamma, \theta)\}^2 d\mathbb{Q}.$$

The quantity  $q_n(\cdot)$  in this criterion is the sample analogue of  $q(\cdot)$  and is called the *functional sample mean squared error* (FSMSE). Under the regularity conditions given above, we can show that the FSMSE converges uniformly to  $q(\cdot)$ . It then follows that the FLS estimator is consistent for  $\theta_*$  and is asymptotically normally distributed around  $\theta_*$ . The following theorem establishes consistency.

**Theorem 2.** *Given Assumptions 1, 2, and 3, as  $n \rightarrow \infty$ ,*

$$(i) \sup_{\theta \in \Theta} |q_n(\theta) - q(\theta)| \rightarrow 0 \text{ a.s.} - \mathbb{P};$$

$$(ii) \hat{\theta}_n \rightarrow \theta_* \text{ a.s.} - \mathbb{P}. \quad \square$$

The uniform consistency of the FSMSE is verified by applying a suitable strong uniform law of large numbers (SULLN). For example, we can apply the SULLN of Andrews (1992) or Newey (1991) under Assumptions 2 and 3. In establishing the SULLN required here we repeatedly invoke the dominated convergence theorem (DCT) to interchange the order of discrete summation and integral operators as  $n \rightarrow \infty$ , for which the moment conditions of Assumption 3 are sufficient. The consistency of the FLS estimator follows directly from the fact that the FSMSE converges to FMSE uniformly on  $\Theta$  a.s.  $-\mathbb{P}$  whenever, as in Assumption 2(iv), the optimizer  $\theta_* := \arg \min_{\theta \in \Theta} q(\theta)$  of the limiting functional  $q(\theta)$  is unique.

Asymptotic normality of the FLS estimator is obtained for this function space setting in parallel to standard derivations for least squares estimation. We begin by observing that by standard Taylor expansion and for some  $\bar{\theta}_n$  between  $\hat{\theta}_n$  and  $\theta_*$ ,

$$\sqrt{n}(\hat{\theta}_n - \theta_*) = A_n^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \int [g_i(\gamma) - \rho_i(\gamma, \theta_*)] \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma), \quad (3)$$

where

$$A_n := \frac{1}{n} \sum_{i=1}^n \int \left\{ \nabla_{\theta} \rho_i(\gamma, \bar{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \bar{\theta}_n) - [g_i(\gamma) - \rho_i(\gamma, \bar{\theta}_n)] \nabla_{\theta}^2 \rho_i(\gamma, \bar{\theta}_n) \right\} d\mathbb{Q}(\gamma).$$

Regular asymptotic behavior in  $A_n$  and  $n^{-1/2} \sum_{i=1}^n \int [g_i(\gamma) - \rho_i(\gamma, \theta_*)] \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma)$  is central to establishing asymptotic normality. The following conditions are sufficient for this purpose.

**Assumption 4. (Hessian Matrix):** *A is positive definite, where  $A := \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \nabla'_{\theta} \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) - \int \int \{\mu(\gamma, x) - \rho(\gamma, \theta_*, x)\} \nabla_{\theta}^2 \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma)$ .*  $\square$

**Assumption 5. (CLT Conditions):** *(i) For each  $j$  and  $j' = 1, 2, \dots, d$ ,  $\int \int \int (\partial/\partial\theta_j) \rho(\gamma, \theta_*, x) \cdot \kappa(\gamma, \tilde{\gamma}|x) \cdot (\partial/\partial\theta_{j'}) \rho(\tilde{\gamma}, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma}) < \infty$ , where  $\kappa(\gamma, \tilde{\gamma}|x) := \int \{g(\gamma) - \rho(\gamma, \theta_*, x)\} \{g(\tilde{\gamma}) - \rho(\tilde{\gamma}, \theta_*, x)\} d\mathbb{P}(g(\gamma), g(\tilde{\gamma})|x)$ ; and*

*(ii) B is positive definite, where  $B := \int \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \kappa(\gamma, \tilde{\gamma}|x) \nabla'_{\theta} \rho(\tilde{\gamma}, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma})$ .*  $\square$

The matrix  $A$  is the probability limit of  $A_n$ , and  $B$  serves as the limiting covariance matrix of the score function in the limit distribution of the FLS estimator. Note that the maximum eigenvalues of  $A$  and  $B$  are both finite by Assumptions 3 and 5(ii),

and the conditional covariance kernel  $\kappa(\cdot, \cdot | x)$  in Assumption 5(i) contributes to the asymptotic covariance matrix of the FLS estimator by way of the matrix  $B$ . Importantly, the presence of the covariance kernel  $\kappa(\cdot, \cdot | x)$  in  $B$  imparts the role of the true probability measure  $\mathbb{P}$  in the limit distribution of the FLS estimator. For a different  $\mathbb{P}$ , a different functional form is obtained for  $\kappa(\cdot, \cdot | x)$ , leading to a different variance matrix  $B$ . In addition,  $B$  also depends on the parametric specification  $\rho_i(\cdot, \theta)$ , implying that different limit distributions are to be expected for different models.

The following theorem establishes asymptotic normality of the FLS estimator, as implied by the regularity conditions imposed so far:

**Theorem 3.** *Given Assumptions 1 – 5,  $\sqrt{n}(\hat{\theta}_n - \theta_*) \stackrel{A}{\sim} N(0, A^{-1}BA^{-1})$ .* □

The result follows in a straightforward manner by applying a strong law and a Lindeberg CLT in conjunction with the Cramér-Wold device to the components of (3), proceeding in parallel to usual derivations for nonlinear least squares estimation as in Newey and McFadden (1994), but executed here within this function space data setting.

The asymptotic covariance matrix exhibits a sandwich form, as is usual. On the other hand, if  $\mathcal{M}$  is correctly specified and if the covariance kernel  $\kappa(\cdot, \cdot | x)$  can be written in functional diagonal form by involving the Dirac delta function, the information matrix equality holds. In particular, note that if  $\mathcal{M}$  is correctly specified,

$$A = \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \nabla'_{\theta} \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma).$$

Further, if we let  $\delta(\cdot)$  be the Dirac delta function and suppose that for some  $\sigma^2 > 0$ ,  $\kappa(\gamma, \gamma' | x) = \sigma^2 \delta(\gamma - \gamma')$  with  $d\mathbb{Q}(\cdot) = \delta(\cdot)$  around 0, it now follows that

$$\begin{aligned} \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \kappa(\gamma, \tilde{\gamma} | x) \nabla'_{\theta} \rho(\tilde{\gamma}, \theta_*, x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma}) &= \sigma^2 \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \delta(\gamma - \tilde{\gamma}) \nabla'_{\theta} \rho(\tilde{\gamma}, \theta_*, x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma}) \\ &= \sigma^2 \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \nabla'_{\theta} \rho(\gamma, \theta_*, x) d\mathbb{Q}(\gamma), \end{aligned}$$

so that  $B = \sigma^2 A$ , implying that the asymptotic covariance matrix reduces to  $\sigma^2 A^{-1}$ . These additional conditions deliver the information equality, thereby motivating model specification testing via the information matrix equality test (e.g., White, 1982; Cho and White, 2014; Cho and Phillips, 2018a).

In general, the information matrix equality holds for the finite dimensional random variable case if the model is correctly specified and the error is conditionally homoskedastic. But these conditions are not sufficient for the information matrix equality in the functional data case. Even if  $\mathcal{M}$  is correct for  $\mu_i(\cdot)$  and  $\varepsilon_i(\cdot) := g_i(\cdot) - \mu_i(\cdot)$  is independent of  $x_i$ , the information matrix equality does not hold. For the information matrix equality, a further condition has to hold for the kernel  $\kappa(\gamma, \tilde{\gamma} | x)$ , whose role is played by Dirac delta function. Note that Dirac delta function is approximable by many popular kernel functions. For example,  $\delta(\gamma - \tilde{\gamma}) = \lim_{s \rightarrow 0} r(\gamma, \tilde{\gamma}; s)$ , where

$$r(\gamma, \tilde{\gamma}; s) := \frac{1}{2\sqrt{s\pi}} \exp\left(-\frac{\|\gamma - \tilde{\gamma}\|^2}{4s}\right), \tag{4}$$

and there are many other kernel functions producing similar results. If  $s$  is sufficiently small,  $B$  can be expected to be approximately proportional to  $A$ .

### Example 1: Distributional Specification – Continued

Given the objective function  $q_n(\cdot)$ , the FLS estimator is now given as follows:

$$\begin{bmatrix} \widehat{\theta}_{1,n} \\ \widehat{\theta}_{2,n} \end{bmatrix} = \begin{bmatrix} \int \rho_1^2(\gamma) d\mathbb{Q}(\gamma) & \int \rho_1(\gamma)\rho_2(\gamma) d\mathbb{Q}(\gamma) \\ \int \rho_1(\gamma)\rho_2(\gamma) d\mathbb{Q}(\gamma) & \int \rho_2^2(\gamma) d\mathbb{Q}(\gamma) \end{bmatrix}^{-1} \begin{bmatrix} n^{-1} \sum_{i=1}^n \int \rho_1(\gamma) g_i(\gamma) d\mathbb{Q}(\gamma) \\ n^{-1} \sum_{i=1}^n \int \rho_2(\gamma) g_i(\gamma) d\mathbb{Q}(\gamma) \end{bmatrix}.$$

The conditions to apply the DCT are applicable to this running example mainly because of the moment conditions in Assumption 3 verified above. These conditions enable application of the SULLN, so that  $\widehat{\mu}_n(\cdot) := n^{-1} \sum_{i=1}^n g_i(\cdot)$  uniformly converges to  $\mu(\cdot)$  with probability 1, and the FLS estimator can be computed using  $\widehat{\mu}_n(\cdot)$ , viz.,

$$\begin{bmatrix} \widehat{\theta}_{1,n} \\ \widehat{\theta}_{2,n} \end{bmatrix} = \begin{bmatrix} \int \rho_1^2(\gamma) d\mathbb{Q}(\gamma) & \int \rho_1(\gamma)\rho_2(\gamma) d\mathbb{Q}(\gamma) \\ \int \rho_1(\gamma)\rho_2(\gamma) d\mathbb{Q}(\gamma) & \int \rho_2^2(\gamma) d\mathbb{Q}(\gamma) \end{bmatrix}^{-1} \begin{bmatrix} \int \rho_1(\gamma) \widehat{\mu}_n(\gamma) d\mathbb{Q}(\gamma) \\ \int \rho_2(\gamma) \widehat{\mu}_n(\gamma) d\mathbb{Q}(\gamma) \end{bmatrix},$$

which now converges to  $\theta_*$  as stated in Theorem 2.

We verify the regularity conditions in Assumptions 4 and 5. The linearity condition of  $\rho(\cdot)$  with respect to  $(\theta_1, \theta_2)$  implies that

$$A = \begin{bmatrix} \int \rho_1^2(\gamma) d\mathbb{Q}(\gamma) & \int \rho_1(\gamma)\rho_2(\gamma) d\mathbb{Q}(\gamma) \\ \int \rho_1(\gamma)\rho_2(\gamma) d\mathbb{Q}(\gamma) & \int \rho_2^2(\gamma) d\mathbb{Q}(\gamma) \end{bmatrix}$$

from the fact that the second-order derivative terms are all zero. As pointed out above,  $A$  is positive definite by Cauchy-Schwarz.

Next,

$$B = \begin{bmatrix} \int \int \rho_1(\gamma) \kappa(\gamma, \widetilde{\gamma}) \rho_1(\widetilde{\gamma}) d\mathbb{Q}(\gamma) d\mathbb{Q}(\widetilde{\gamma}) & \int \int \rho_1(\gamma) \kappa(\gamma, \widetilde{\gamma}) \rho_2(\widetilde{\gamma}) d\mathbb{Q}(\gamma) d\mathbb{Q}(\widetilde{\gamma}) \\ \int \int \rho_2(\gamma) \kappa(\gamma, \widetilde{\gamma}) \rho_1(\widetilde{\gamma}) d\mathbb{Q}(\gamma) d\mathbb{Q}(\widetilde{\gamma}) & \int \int \rho_2(\gamma) \kappa(\gamma, \widetilde{\gamma}) \rho_2(\widetilde{\gamma}) d\mathbb{Q}(\gamma) d\mathbb{Q}(\widetilde{\gamma}) \end{bmatrix},$$

where  $\kappa(\gamma, \widetilde{\gamma}) := \mathbb{E}[(g_i(\gamma) - \rho(\gamma, \theta_{1*}, \theta_{2*}))(g_i(\widetilde{\gamma}) - \rho(\widetilde{\gamma}, \theta_{1*}, \theta_{2*}))]$ , which is required to be positive definite by Assumption 5.

In general, it is not straightforward to verify this condition, but it can be confirmed by estimating  $B$  consistently as discussed in Section 3.3. If the mixture assumption is correct with  $\pi_{\dagger} = 0$ ,  $B$  can be numerically computed by noting that

$$\kappa(\gamma, \widetilde{\gamma}) = \mathbb{E}[g_i(\gamma)g_i(\widetilde{\gamma})] = \frac{(2\gamma - 1)^{1/2}(2\widetilde{\gamma} - 1)^{1/2}}{\gamma + \widetilde{\gamma} - 1}. \quad \square$$

## 3.2 FLS Estimation with Nuisance Effects

Functional data analysis often involves nuisance effects by the very nature of the data. In particular, when data are constructed using aligned discrete observations, such a construction naturally introduces nuisance effects. To examine this extension, we characterize the functional data in the form of the transform

$$\widetilde{g}_i : \Gamma \times \Xi \mapsto \mathbb{R},$$

where  $\Gamma$  is the same as before, and  $\Xi$  is a compact parameter space for a nuisance parameter  $\xi_*$ , so that the functional observations are defined on the product space  $\Gamma \times \Xi$ . We assume that the nuisance parameter  $\xi_*$  is identifiable and can be consistently estimated by  $\widehat{\xi}_n$  obtained in a preliminary stage before estimating  $\theta_*$ , from which our functional observations are constructed as  $\widehat{g}_i(\cdot) \equiv \widetilde{g}_i(\cdot, \widehat{\xi}_n)$ . This assumption on the data structure generalizes that assumed in Section 3.1 because for some known  $\xi_*$ , we can let  $g_i(\cdot)$  be identical to  $\widetilde{g}_i(\cdot, \xi_*)$ . So the data analysis given in this section is also applicable to  $\{g_i(\cdot), x_i\}_{i=1}^n$ .

Notwithstanding this specialization, the asymptotic influence of the nuisance effects on the FLS estimator is not negligible in general and typically modifies the limit behavior of the FLS estimator. The results of Section 3.1 are therefore extended here to accommodate the effects conveyed by nuisance parameter estimation.

Functional data are often influenced by nuisance effects, as described in the Introduction. First, when a model is unidentified under the null (e.g., Davies, 1977, 1987), a functional data set with nuisance effects can be collected by letting each individual observation be the score function defined on the set of unidentified parameters with the other parameters being evaluated at the parameter estimates obtained using the null model. In such cases the null parameter estimates play the role of  $\widehat{\xi}_n$ , and  $\gamma$  can be treated as generic notation for the unidentified parameters with  $\widetilde{g}_i(\cdot, \cdot)$  being the score function. Second, functional observations are often constructed by local polynomial kernel, sieve, kriging, and polynomial spline estimation using discrete observations (e.g., Zhang and Chen, 2007; Chen, 2007; Chilès and Delfiner, 1999; Wang, 2011). In these cases functional observations are influenced by the kernel or sieve estimation error that is captured by  $\widehat{\xi}_n$  here, thereby potentially modifying the large sample properties of the FLS estimator. Specifically, using an optimal bandwidth choice, Zhang and Chen (2007) show that estimated functional observations obtained by local polynomial kernel estimation uniformly converge to continuous functional observations at the rate  $n^{(p+1)\delta/(2p+3)}$ , where  $p$  is the degree of polynomial function and  $\delta$  is a positive number such that  $n_t \geq Cn^\delta$  uniformly in  $t$  for some positive number  $C$ , and  $n_t$  is the number of discrete observations underlying the  $t$ -th functional observation. If a sufficiently large  $p$  is selected such that  $(p+1)\delta/(2p+3) > 1/2$ , the functional approximation obtained by local polynomial kernel estimation is super-consistent for the latent functional observation, so that the nuisance effect in the approximated functional observation can be ignored in the limit theory when applying FLS estimation because the convergence rate of the FLS estimator is  $\sqrt{n}$ , as given in Theorem 3. As discussed in Section 5.1 using income processes, if  $p$  is sufficiently large, the key condition for the use of the optimal bandwidth can be satisfied even if  $\delta$  is not so different from unity, implying that the estimated functional observations (with  $n_t$  proportional to  $n$ ) can be approximated by those estimated with the optimal bandwidth. On the other hand, if the key condition does not hold because  $\delta < 1$ , we cannot use the Zhang and Chen (2007) optimal bandwidth. Nevertheless, the functional observations can still be consistently estimated from discrete observations by letting the bandwidth converge to zero at a rate slower than  $\delta$ .

To fix ideas let the data set be given as  $\{(\widehat{g}_i(\cdot), x_i')\}_{i=1}^n$ . After replacing  $g_i$  with  $\widehat{g}_i$ , we obtain the FLS estimator by minimizing the functional mean squared error as before, viz.,

$$\widetilde{\theta}_n := \arg \min_{\theta \in \Theta} \widehat{q}_n(\theta), \quad \text{where} \quad \widehat{q}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \int \{\widehat{g}_i(\gamma) - \rho_i(\gamma, \theta)\}^2 d\mathbb{Q}.$$

Henceforth, we refer to  $\widetilde{\theta}_n$  as the two-stage FLS (TSFLS) estimator for parametric estimation in the mean function.

We now proceed to examine how the estimation error imbedded in  $\widehat{\xi}_n$  changes the asymptotic behavior of the FLS estimator. To tackle this issue we start by extending the previous regularity conditions for Theorems 2 and 3 to cope with the presence of nuisance effects. We modify Assumptions 1 and 3 in the following:

**Assumption 6. (Data):** (i) Let  $(\Omega, \mathcal{F}, \mathbb{P})$  and  $(\Gamma, \tau)$  be a complete probability space and a compact metric space respectively.  $\Gamma \subset \mathbb{R}^d$  ( $d \in \mathbb{N}$ ) and  $\Xi \subset \mathbb{R}^s$  ( $s \in \mathbb{N}$ ) are compact;

(ii)  $\{(\widetilde{g}_i(\cdot), x_i') : \widetilde{g}_i : \Omega \times \Gamma \times \Xi \mapsto \mathbb{R} \text{ and } x_i : \Omega \mapsto \mathbb{R}^k\}_{i=1}^n$  ( $k \in \mathbb{N}$ ) is a set of IID observations such that

(ii.a) for each  $(\gamma, \xi) \in \Gamma \times \Xi$ ,  $(\widetilde{g}_i(\cdot, \gamma, \xi), x_i')$  is measurable –  $\mathcal{F}$ ;

(ii.b) for each  $\xi \in \Xi$ ,  $\widetilde{g}_i(\omega, \cdot, \xi) \in \mathcal{C}^{(0)}(\Gamma)$  for all  $\omega \in F \in \mathcal{F}$  with  $\mathbb{P}(F) = 1$ ;

(ii.c) for each  $(\omega, \gamma) \in \Omega \times \Gamma$ ,  $\widetilde{g}_i(\omega, \gamma, \cdot)$  is in  $\mathcal{C}^{(1)}(\Xi)$  for all  $\omega \in F \in \mathcal{F}$  with  $\mathbb{P}(F) = 1$ ;

(iii)  $(\Gamma, \mathcal{G}, \mathbb{Q})$  and  $(\Omega \times \Gamma, \mathcal{F} \otimes \mathcal{G}, \mathbb{P} \times \mathbb{Q})$  are complete probability spaces and for  $i = 1, 2, \dots$  and  $\xi \in \Xi$ ,  $\tilde{g}_i(\cdot, \cdot, \xi)$  is measurable  $-\mathcal{F} \otimes \mathcal{G}$ .  $\square$

**Assumption 7. (E-Moments):** For some  $m_i \in L^2(\mathbb{P})$ ,

$$(i) \sup_{(\gamma, \xi) \in \Gamma \times \Xi} |\tilde{g}_i(\gamma, \xi)| \leq m_i \text{ a.s. } -\mathbb{P};$$

$$(ii) \sup_{(\gamma, \theta) \in \Gamma \times \Theta} |\rho_i(\gamma, \theta)| \leq m_i \text{ a.s. } -\mathbb{P};$$

$$(iii) \sup_j \sup_{(\gamma, \xi) \in \Gamma \times \Xi} |(\partial/\partial \xi_j) \tilde{g}_i(\gamma, \xi)| \leq m_i \text{ a.s. } -\mathbb{P};$$

$$(iv) \text{ for each } j = 1, 2, \dots, d, \sup_{(\gamma, \theta, \xi) \in \Gamma \times \Theta \times \Xi} |(\partial/\partial \theta_j) \rho_i(\gamma, \theta, \xi)| \leq m_i \text{ a.s. } -\mathbb{P};$$

$$(v) \text{ for each } j \text{ and } j' = 1, 2, \dots, d, \sup_{(\gamma, \theta, \xi) \in \Gamma \times \Theta \times \Xi} |(\partial^2/\partial \theta_j \partial \theta_{j'}) \rho_i(\gamma, \theta, \xi)| \leq m_i \text{ a.s. } -\mathbb{P}. \quad \square$$

Consistency of the TSFLS estimator can be verified by investigating the limit behavior of the first-order conditions for the TSFLS estimator. For this purpose, note that for some  $\bar{\xi}_{n, \gamma}$  between  $\hat{\xi}_n$  and  $\xi_*$ , we have

$$\begin{aligned} & \frac{1}{n} \int \sum_{i=1}^n \{\hat{g}_i(\gamma) - \rho_i(\gamma, \theta_*)\} \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma) \\ &= \frac{1}{n} \int \sum_{i=1}^n \{\tilde{g}_i(\gamma, \xi_*) - \rho_i(\gamma, \theta_*)\} \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma) + \frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \theta_*) [\nabla'_{\xi} \tilde{g}_i(\gamma, \bar{\xi}_{n, \gamma})] d\mathbb{Q} \cdot (\hat{\xi}_n - \xi_*). \end{aligned} \quad (5)$$

The left side of (5) is the first-order derivative of  $\hat{q}_n(\cdot)$  with respect to  $\theta$  evaluated at  $\theta_*$ , whereas the right side is the Taylor expansion with respect to  $\xi$  around  $\xi_*$ . Assumption 7(iii) enables use of the mean-value theorem. Proving consistency of  $\tilde{\theta}_n$  for  $\theta_*$  then involves showing that the quantities on the right side of (5) vanish as  $n \rightarrow \infty$ . Since the first component in (5) trivially vanishes by applying the proof of Theorem 3, the result involves showing that the second term converges to zero in probability. Note that the second component is  $O_{\mathbb{P}}(\hat{\xi}_n - \xi_*)$  because the sample average of the integrals in the second term is  $O_{\mathbb{P}}(1)$  by Assumption 7, so that if the deviation  $(\hat{\xi}_n - \xi_*)$  is asymptotically negligible, the first-order condition asymptotically holds at  $\theta_*$ , leading to consistency of the TSFLS estimator. For this purpose we impose the following high level condition concerning  $\hat{\xi}_n$ .

**Assumption 8. (E-Estimator-1):** There exists a sequence of measurable functions  $\{\hat{\xi}_n : \Omega \mapsto \Xi\}$  such that

$$(i) \hat{\xi}_n \rightarrow \xi_* \text{ a.s. } -\mathbb{P}, \text{ where } \xi_* \text{ is an interior element in } \Xi. \quad \square$$

Almost sure convergence of  $\hat{\xi}_n$  is more restrictive here than convergence in probability and could be relaxed, but is used because the condition assists in simplifying subsequent derivations. Consistency of the TSFLS estimator is then immediate under these conditions.

**Theorem 4.** Given Assumptions 2, 6, 7, and 8(i), as  $n$  tends to infinity,

$$(i) \sup_{\theta \in \Theta} |\hat{q}_n(\theta) - q(\theta)| \rightarrow 0 \text{ a.s. } -\mathbb{P}; \text{ and}$$

$$(ii) \tilde{\theta}_n \rightarrow \theta_* \text{ a.s. } -\mathbb{P}. \quad \square$$

The assumptions for Theorem 3 are insufficient to deliver the limit distribution of the TSFLS estimator as they do not address the asymptotic properties of  $\hat{\xi}_n$ . To establish asymptotic normality we impose the following additional conditions.

**Assumption 8. (E-Estimator-2):** (ii) there exists a finite nonstochastic  $s \times s$  positive definite matrix  $H$  and a sequence of random vectors  $\{s_{*n}\}$  measurable  $-\mathcal{F}$  for which  $\sqrt{n}(\hat{\xi}_n - \xi_*) = -H^{-1} \sqrt{n} s_{*n} + o_{\mathbb{P}}(1)$ ;

$$(iii) \text{ for } i = 1, 2, \dots, \text{ there exists } s_i : \Omega \times \Xi \mapsto \mathbb{R}^s \text{ such that:}$$

$$(iii.a) \text{ for each } \xi \in \Xi, s_i(\cdot, \xi) \text{ is measurable } -\mathcal{F} \text{ and IID over } i;$$

(iii.b)  $s_i(\omega, \cdot)$  is continuous for all  $\omega \in F \subset \mathcal{F}$ ,  $\mathbb{P}(F) = 1$ ;

(iii.c) for some  $m_i \in L^2(\mathbb{P})$ ,  $|s_i(\omega, \cdot)| \leq m_i(\omega)$ ; and

(iii.d)  $\sqrt{n}s_{*n} = n^{-1/2} \sum_{i=1}^n s_i(\cdot, \xi_*) + o_{\mathbb{P}}(1)$  such that for each  $j = 1, 2, \dots, s$ ,  $\mathbb{E}_{\mathbb{P}}[s_{ji}(\cdot, \xi_*)^2] < \infty$ , where  $s_{ji}(\cdot, \xi_*)$  is the  $j$ -th row element of  $s_i(\cdot, \xi_*)$ .  $\square$

Assumptions 8(ii and iii) assume that  $\widehat{\xi}_n - \xi_*$  is asymptotically equivalent to the product of the nonstochastic matrix  $H$  and the score  $s_{*n}$ . Many estimators satisfy this characteristic asymptotically, including least squares, generalized method of moments, and (quasi-)maximum likelihood estimators. In order to retain generality in the analysis, we do not specify here how  $H$  and  $s_{*n}$  are obtained from primitive model formulations. Treating  $s_{*n}$  as being formally defined in Assumption 8, we now use  $s_i(\xi)$  to denote  $s_i(\omega, \xi)$ , suppressing the argument  $\omega$  for notational ease.

**Assumption 9. (E-CLT):** Let  $J := \mathbb{E}[s_i(\xi_*)s_i(\xi_*)']$ ,  $K := \int \mathbb{E}_{\mathbb{P}}[s_i(\xi_*)\{\tilde{g}_i(\gamma, \xi_*) - \rho_i(\gamma, \theta_*)\}\nabla'_{\theta}\rho_i(\gamma, \theta_*)]d\mathbb{Q}(\gamma)$ , and  $B$  be as defined earlier in Assumption 5. Let

$$C := \begin{bmatrix} J & K \\ K' & B \end{bmatrix},$$

and assume the following:

(i)  $C$  is positive definite;

(ii)  $B_*$  is positive definite, where  $B_* := B - MH^{-1}K - K'H^{-1}M' + MH^{-1}JH^{-1}M'$  and  $M := \int \mathbb{E}_{\mathbb{P}}[\nabla_{\theta}\rho_i(\gamma, \theta_*)\nabla'_{\xi}\tilde{g}_i(\gamma, \xi_*)]d\mathbb{Q}(\gamma)$ .  $\square$

Assumption 9 characterizes the key components needed for the limiting covariance matrix of the TSFLS estimator in the presence of nuisance effects. It generalizes Assumption 5 to accommodate the additional estimator  $\widehat{\xi}_n$ . The matrix  $C$  is employed to capture the asymptotic covariance matrix between the score vectors of the estimates  $\widehat{\xi}_n$  and  $\tilde{\theta}_n$ , thereby providing a channel for the nuisance effects to be conveyed to the limit distribution of the TSFLS estimator. With this framework for the nuisance effects in hand, asymptotic normality of the TSFLS estimator is established in the following theorem.

**Theorem 5.** Given Assumptions 2, 4, 6, 7, 8, and 9,  $\sqrt{n}(\tilde{\theta}_n - \theta_*) \overset{\Delta}{\sim} \mathcal{N}(0, A^{-1}B_*A^{-1})$ .  $\square$

Just as for Theorem 3 without nuisance parameters, the asymptotic distribution of the TSFLS estimator is obtained in a standard fashion (Newey and McFadden, 1994). If  $\xi_*$  were known and we can let  $s_i(\xi_*) \equiv 0$ , then  $B_*$  reduces to  $B$  by definition as  $K$  and  $J$  are a zero vector and matrix, respectively. But as the result makes clear, if the nuisance effect is not asymptotically negligible, the asymptotic variance matrix of  $\tilde{\theta}_n$  changes from  $A^{-1}BA^{-1}$  to  $A^{-1}B_*A^{-1}$ , thereby modifying the limit variability of  $\tilde{\theta}_n$ . So the covariance matrix is typically adjusted when the parameter of interest is estimated via use of another parameter estimate. Among others, for example, Amemiya (1979) and Lee, Maddala, and Trost (1980) examine estimating a structural parameter by a two-stage method in a limited dependent variable context, and the resulting variance matrix of their estimator is affected by the first-stage estimator, in a similar way as  $B_*$ . The presence of nuisance effects introduces further changes in the construction of appropriate test statistic formulae, as shown below.

Before proceeding two further remarks are in order. First, estimating the unconditional mean function of functional data can be conducted in parallel to the estimation of the conditional mean function. In view of this similarity we leave that discussion to the Appendix. Second, the asymptotic efficiency of the FLS estimator depends on the adjunct probability measure  $\mathbb{Q}$ . Theorem 1 implies that the argument minimizing  $q(\cdot)$  is affected by  $\mathbb{Q}$  if  $\mathcal{M}$  is misspecified. In addition, the asymptotic covariance matrices in Theorem 3 and 5 have different values depending on  $\mathbb{Q}$ . This property implies that an asymptotically more efficient FLS

estimator may be obtained by parameterizing the adjunct probability measure to nest the probability measure of interest as a special case of the parameterized adjunct probability measure. For example, if  $\Gamma \subset \mathbb{R}$ , a beta distribution could be assumed for  $\mathbb{Q}$ , nesting a uniform distribution by setting the shape parameters to unity. With this extension the FLS estimator can be obtained by minimizing the FMSE incorporating the beta distribution with respect to both the model parameters and the shape parameters, thereby optimizing the adjunct probability measure in addition to the FMSE. The asymptotic efficiency of this extended FLS estimator can be compared with the FLS estimator obtained by when  $\mathbb{Q}$  is uniform. In the event that the optimized adjunct probability measure substantially differs from the uniform distribution, an asymptotic efficiency gain might be expected from the extended FLS estimator.

### 3.3 Estimation of the Covariance Matrix of FLS

The role of the covariance matrices in Theorems 3 and 5 is important as these matrices are used to construct test statistics. This section examines how these covariance matrices may be estimated consistently.

First, we discuss the case with no nuisance effects where the covariance matrix is  $A^{-1}BA^{-1}$  in Theorem 3. The domination conditions in Assumption 3 enable application of the SULLN to the estimators given by  $\widehat{A}_n$  and  $\widehat{B}_n$  in Theorem 6 below, so that  $\widehat{A}_n^{-1}\widehat{B}_n\widehat{A}_n^{-1}$  provides a consistent estimator of the covariance matrix.

**Theorem 6.** *Given Assumptions 1, 2, 3, and 5,  $\widehat{A}_n \rightarrow A$  a.s.  $-\mathbb{P}$  and  $\widehat{B}_n \rightarrow B$  a.s.  $-\mathbb{P}$ , where*

$$\begin{aligned}\widehat{A}_n &:= \frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \widehat{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \widehat{\theta}_n) d\mathbb{Q}(\gamma) - \frac{1}{n} \sum_{i=1}^n \int \varepsilon_i(\gamma, \widehat{\theta}_n) \nabla_{\theta}^2 \rho_i(\gamma, \widehat{\theta}_n) d\mathbb{Q}(\gamma), \\ \widehat{B}_n &:= \frac{1}{n} \sum_{i=1}^n \int \int \nabla_{\theta} \rho_i(\gamma, \widehat{\theta}_n) \varepsilon_i(\gamma, \widehat{\theta}_n) \varepsilon_i(\widetilde{\gamma}, \widehat{\theta}_n) \nabla'_{\theta} \rho_i(\widetilde{\gamma}, \widehat{\theta}_n) d\mathbb{Q}(\gamma) d\mathbb{Q}(\widetilde{\gamma}),\end{aligned}$$

and for each  $\gamma$  and  $\theta$ ,  $\varepsilon_i(\gamma, \theta) := g_i(\gamma) - \rho_i(\gamma, \theta)$ . □

Next consider models with nuisance effects. Since  $B_*$  involves component matrices in its definition, further conditions are needed to ensure consistent estimation. The following conditions are used for this purpose.

**Assumption 10. (E-Covariance):** (i) *For a sequence of measurable functions  $\{\widehat{J}_n : \Omega \mapsto \mathbb{R}^{s \times s}\}$ ,  $\widehat{J}_n \rightarrow J$  a.s.  $-\mathbb{P}$ ; and*  
(ii) *for a sequence of measurable functions  $\{\widehat{H}_n : \Omega \mapsto \mathbb{R}^{s \times s}\}$ ,  $\widehat{H}_n \rightarrow H$  a.s.  $-\mathbb{P}$ .* □

**Assumption 11. (SULLN):** *Let  $\varepsilon_i(\gamma, \theta, \xi) := g_i(\gamma, \xi) - \rho_i(\gamma, \theta)$ . For some  $m_i \in L^2(\mathbb{P})$ ,*

- (i)  $\sup_{(\gamma, \theta, \xi) \in \Gamma \times \Theta \times \Xi} |\varepsilon_i(\gamma, \theta, \xi)| \leq m_i$  a.s.  $-\mathbb{P}$ ;
- (ii) for  $j = 1, 2, \dots, d$ ,  $\sup_{(\gamma, \theta) \in \Gamma \times \Theta} |(\partial/\partial\theta_j)\rho_i(\gamma, \theta)| \leq m_i$  a.s.  $-\mathbb{P}$ ;
- (iii) for each  $j$  and  $j' = 1, 2, \dots, d$ ,  $\sup_{(\gamma, \theta, \xi) \in \Gamma \times \Theta \times \Xi} |(\partial^2/\partial\theta_j\partial\theta_{j'})\rho_i(\gamma, \theta)| \leq m_i$  a.s.  $-\mathbb{P}$ ;
- (iv) for each  $j = 1, 2, \dots, s$ ,  $\sup_{(\gamma, \xi) \in \Gamma \times \Xi} |(\partial/\partial\xi_j)\widetilde{g}_i(\gamma, \xi)| \leq m_i$  a.s.  $-\mathbb{P}$ ; and
- (v) for each  $j = 1, 2, \dots, s$ ,  $\sup_{\xi \in \Xi} |s_{ji}(\xi)| \leq m_i$  a.s.  $-\mathbb{P}$ , where  $s_{ji}(\xi)$  is the  $j$ -th row element of  $s_i(\xi)$ . □

Under Assumption 10, the two submatrices  $H$  and  $J$  appearing in  $B_*$  can be consistently estimated by using  $\widehat{H}_n$  and  $\widehat{J}_n$ . In general, these estimators are obtained by preliminary estimation of  $\widehat{\xi}_n$  and can be easily computed using standard methods. For example, if  $\widehat{\xi}_n$  is a (quasi-) maximum likelihood estimator,  $\widehat{H}_n$  and  $\widehat{J}_n$  may be identified as the Hessian matrix of the quasi-likelihood function and the sample average of the products of the first-order derivatives evaluated at  $\widehat{\xi}_n$  in the usual fashion. Note that Assumptions 11(ii and iii) are stronger than Assumption 3 because the SULLN is required to hold not only for the parameter



space  $\Gamma \times \Theta$  but for  $\Xi$  as well. Furthermore, Assumptions 11(iii and iv) require the SULLN to hold for other random elements used to provide consistent estimation of the component matrices  $K$  and  $M$  of  $B_*$  that appear in Assumption 9.

The following Theorem provides consistent estimators for  $A$  and  $B_*$  under these regularity conditions.

**Theorem 7.** Let  $\tilde{\varepsilon}_{in}(\gamma, \theta) := \varepsilon(\gamma, \theta, \hat{\xi}_n)$ . Given Assumptions 2, 6, 8, 9, 10, and 11,  $\tilde{A}_n \rightarrow A$  a.s. -  $\mathbb{P}$  and  $\tilde{B}_n \rightarrow B_*$  a.s. -  $\mathbb{P}$ , where

$$\begin{aligned}\tilde{A}_n &:= \frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \tilde{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \tilde{\theta}_n) d\mathbb{Q}(\gamma) - \frac{1}{n} \sum_{i=1}^n \int \tilde{\varepsilon}_{in}(\gamma, \tilde{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \tilde{\theta}_n) d\mathbb{Q}(\gamma); \\ \tilde{B}_n &:= \bar{B}_n - \widehat{M}_n \widehat{H}_n^{-1} \widehat{K}_n - \widehat{K}'_n \widehat{H}_n^{-1'} \widehat{M}'_n + \widehat{M}_n \widehat{H}_n^{-1} \widehat{J}_n \widehat{H}_n^{-1'} \widehat{M}'_n; \\ \bar{B}_n &:= \frac{1}{n} \sum_{i=1}^n \int \int \nabla_{\theta} \rho(\gamma, \tilde{\theta}_n) \tilde{\varepsilon}_{in}(\gamma, \tilde{\theta}_n) \tilde{\varepsilon}_{in}(\tilde{\gamma}, \tilde{\theta}_n) \nabla'_{\theta} \rho(\tilde{\gamma}, \tilde{\theta}_n) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma}); \\ \widehat{M}_n &:= \frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \tilde{\theta}_n) \nabla'_{\xi} \tilde{g}_i(\gamma, \hat{\xi}_n) d\mathbb{Q}(\gamma); \text{ and } \widehat{K}_n := \frac{1}{n} \sum_{i=1}^n \int s_i(\tilde{\theta}_n) \tilde{\varepsilon}_{in}(\gamma, \tilde{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \tilde{\theta}_n) d\mathbb{Q}(\gamma). \quad \square\end{aligned}$$

From these results it follows that the covariance matrix in Theorem 5 can be consistently estimated by  $\tilde{A}_n^{-1} \tilde{B}_n \tilde{A}_n^{-1}$ . If  $M = 0$ , the limit distribution of the TSFLS estimator is identical to that of FLS estimator. Hence, both estimators  $\tilde{A}_n^{-1} \tilde{B}_n \tilde{A}_n^{-1}$  and  $\tilde{A}_n^{-1} \bar{B}_n \tilde{A}_n^{-1}$  are consistent.

## Example 2: Inference on Random Coefficients

Standard regression models typically assume the coefficients of explanatory variables are fixed parameters. When such assumptions are violated, statistical inference can be misleading due to biases in estimating the standard errors. Inference using the random coefficient model is an alternative approach that is particularly useful in modeling conditional heteroskedasticity or time varying coefficients in time series models. Many studies in the literature consider testing for random coefficients in regression models (e.g., see Hsiao, 1974; Breusch and Pagan, 1979; Ramanathan and Rajarshi, 1992; Swamy and Tavlas, 1995; Akharif, Fihri, Hallin, and Mellouk 2020).

Consider the simple linear regression model

$$y_i = x_i' \beta_i + \delta_i^{1/2} \varepsilon_i, \quad (6)$$

with  $x_i := (1, z_i)'$  and  $z_i \in \mathbb{R}$  an explanatory variable. The regression coefficient  $\beta_i$  is formulated to embody a potential random element so that

$$\beta_i := (\psi_{1\ddagger}, \psi_{2\ddagger})' + \pi_{\ddagger}^{1/2} \Omega^{1/2}(\gamma_{\ddagger}) \nu_i. \quad (7)$$

where  $(\psi_{1\ddagger}, \psi_{2\ddagger})'$  is a constant vector and  $\nu_i \in \mathbb{R}^2$  is a random vector with variance  $I_2$ . The matrix coefficient function  $\Omega(\gamma)$  is assumed to be positive definite uniformly on the space  $\Gamma$  to which the true parameter  $\gamma_{\ddagger}$  belongs, and  $\pi_{\ddagger} \geq 0$ . The coefficient  $\beta_i$  in (7) is constant if and only if  $\pi_{\ddagger} = 0$ .

This type of random coefficient model is commonly studied and used in the applied literature. For example, Andrews (2001) developed limit theory for models similar to the random coefficient model where parameters may lie on the boundary of the parameter space, such as  $\pi_{\ddagger} = 0$  in (7). Rosenberg (1973) and Engle and Watson (1985) extended the random coefficient model to include conditional heteroskedastic processes in time series settings; and many empirical studies exploited features of the random coefficient model relevant to investigating conditional heterogeneity in time series data (e.g., Swamy and Tinsley, 1980; Stock and Watson, 1998).

In what follows we relate the random coefficient model to FDA model analysis. We first note that substituting the expression of  $\beta_i$  into (6) yields the following conditional heteroskedasticity model  $y_i = x_i' \psi_{\dagger} + u_i$ , where  $\psi_{\dagger} := (\psi_{1\dagger}, \psi_{2\dagger})'$  and  $u_i := \pi_{\dagger}^{1/2} x_i' \Omega^{1/2}(\gamma_{\dagger}) \nu_i + \delta_{\dagger}^{1/2} \varepsilon_i$ , which leads to the explicit conditional variance function

$$\text{var}(u_i | x_i) = \delta_{\dagger} + \pi_{\dagger} x_i' \Omega(\gamma_{\dagger}) x_i = \delta_{\dagger} + \pi_{\dagger} [1 + \exp(\gamma_{\dagger}) z_i^2] \quad (8)$$

when the variance matrix function  $\Omega(\gamma_{\dagger})$  has the form

$$\Omega(\gamma_{\dagger}) := \begin{bmatrix} 1 & 0 \\ 0 & \exp(\gamma_{\dagger}) \end{bmatrix}.$$

Let  $\Gamma = \{\gamma : \gamma \in [0, 1]\}$  and suppose that the researcher estimates the unknown parameter values by maximum likelihood under Gaussian assumptions, so that the data are generated according to  $(\varepsilon_i, \nu_i)' | z_i \sim \text{IID } \mathcal{N}(0, I_3)$  and  $z_i \sim \text{IID } U[0, 1]$ , implying that the random coefficient  $\beta_i$  is assumed to follow a normal distribution. This assumption is a standard one, as assumed by Rosenberg in a Bayesian perspective (1973).

In spite of this detailed specification for maximum likelihood estimation of the unknown parameters, it is not straightforward to test  $\pi_{\dagger} = 0$ . If  $\pi_{\dagger} = 0$  then  $\gamma_{\dagger}$  in (8) is not identified, as in the earlier example, so that the likelihood ratio test statistic (say) does not have a chi-squared limit distribution under the null in view of the Davies (1977, 1987) identification problem unless  $\gamma_{\dagger}$  is fixed at a particular value. This difficulty means that additional effort is needed to obtain critical values for a test such as the  $C(\alpha)$  test statistic detailed in Section 5.2.

We propose to test the random coefficient property within the FDA model setting. For this purpose, we reformulate the given DGP in the FLS framework in parallel to the earlier example. Then, the score with respect to  $\pi$  at  $\pi_{\dagger} = 0$  (i.e., non-random coefficient) is obtained as

$$\frac{1}{2\delta_{\dagger}^2} \sum_{i=1}^n [1 + \exp(\gamma_{\dagger}) z_i^2] \{(y_i - x_i' \psi_{\dagger})^2 - \delta_{\dagger}\}, \quad (9)$$

and  $\mathbb{E}[(y_i - x_i' \psi_{\dagger})^2 | z_i] = \delta_{\dagger}$ , so that the conditional mean of (9) is zero irrespective of  $\gamma \in \Gamma = [0, 1]$ . On the other hand, if the coefficient is random, the population mean of (9) is obtained as  $\frac{n\pi_{\dagger}}{2\delta_{\dagger}^2} \mathbb{E}\{[1 + \exp(\gamma_{\dagger}) z_i^2][1 + \exp(\gamma_{\dagger}) z_i^2]\}$  using (8), thereby motivating the following random function as a candidate function for  $\tilde{g}_i$ :

$$\tilde{g}_i(\gamma, \psi, \delta) := \{1 + \exp(\gamma) z_i^2\} \{(y_i - x_i' \psi)^2 - \delta\},$$

where we can estimate the unknown parameters  $\psi_{\dagger}$  and  $\delta_{\dagger}$  by the quasi-maximum likelihood estimator assuming a normal distribution for the error term, viz.,  $\hat{\psi}_n := (\sum_{i=1}^n x_i x_i')^{-1} \sum_{i=1}^n x_i y_i$  and  $\hat{\delta}_n := \frac{1}{n} \sum_{i=1}^n (y_i - x_i' \hat{\psi}_n)^2$ , respectively, which are consistent for  $\xi_* := (\psi_*', \delta_*)' := (\psi_{\dagger}', \delta_{\dagger} + \pi_{\dagger} [1 + \exp(\gamma_{\dagger}) \mathbb{E}[z_i^2]])'$  and  $\mathbb{E}[z_i^2] = \frac{1}{3}$  for this particular example. That is, the quasi-maximum likelihood estimators are able to estimate the coefficient in the conditional mean and the unconditional variance of the error. Note that if  $\pi_{\dagger} = 0$ , then  $\delta_* = \delta_{\dagger}$ . Accordingly, our functional observations can be constructed as  $\hat{g}_i(\gamma) := \tilde{g}_i(\gamma, \hat{\delta}_n, \hat{\psi}_n)$ , and for this, we may specify a model for the mean function as follows:

$$\rho_i(\gamma, \theta) := \theta_1 \rho_{1i}(\gamma) + \theta_2 \rho_{2i}(\gamma) := \theta_1 (1 + \exp(\gamma) z_i^2) + \theta_2 (z_i^2 + \exp(\gamma) z_i^4)$$

by noting that for each  $\gamma$ ,  $\mu_i(\gamma) = \mathbb{E}[\tilde{g}_i(\gamma, \psi_*, \delta_*) | z_i] = -\pi_{\dagger} \exp(\gamma_{\dagger}) (1 + \exp(\gamma) z_i^2) + \pi_{\dagger} \exp(\gamma_{\dagger}) (z_i^2 + \exp(\gamma) z_i^4)$ , implying

that  $\theta_* = (\theta_{1*}, \theta_{2*})' = (-\pi_{\dagger} \exp(\gamma_{\dagger}), \pi_{\dagger} \exp(\gamma_{\dagger}))'$ . That is, this model is correctly specified by letting  $\rho_{1i}(\cdot)$  and  $\rho_{2i}(\cdot)$  be  $(1 + \exp(\cdot)z_i^2)$  and  $(z_i^2 + \exp(\cdot)z_i^4)$ , respectively. Note that  $(\theta_{1*}, \theta_{2*})$  is also obtained by minimizing  $q(\cdot)$ , viz.,

$$\theta_* = \begin{bmatrix} \int \int \rho_{1i}^2(\gamma) d\mathbb{P}(z) d\mathbb{Q}(\gamma) & \int \int \rho_{1i}(\gamma) \rho_{2i}(\gamma) d\mathbb{P}(z) d\mathbb{Q}(\gamma) \\ \int \int \rho_{1i}(\gamma) \rho_{2i}(\gamma) d\mathbb{P}(z) d\mathbb{Q}(\gamma) & \int \int \rho_{2i}^2(\gamma) d\mathbb{P}(z) d\mathbb{Q}(\gamma) \end{bmatrix}^{-1} \begin{bmatrix} \int \int \rho_{1i}(\gamma) \mu_i(\gamma) d\mathbb{P}(z) d\mathbb{Q}(\gamma) \\ \int \int \rho_{2i}(\gamma) \mu_i(\gamma) d\mathbb{P}(z) d\mathbb{Q}(\gamma) \end{bmatrix},$$

where the integrals are taken over  $\Gamma = [0, 1]$ . The TSFLS estimator is now obtained as a sample analog to  $\theta_*$ , viz.,

$$\hat{\theta}_n = \begin{bmatrix} \sum_{i=1}^n \int \rho_{1i}^2(\gamma) d\mathbb{Q}(\gamma) & \sum_{i=1}^n \int \rho_{1i}(\gamma) \rho_{2i}(\gamma) d\mathbb{Q}(\gamma) \\ \sum_{i=1}^n \int \rho_{1i}(\gamma) \rho_{2i}(\gamma) d\mathbb{Q}(\gamma) & \sum_{i=1}^n \int \rho_{2i}^2(\gamma) d\mathbb{Q}(\gamma) \end{bmatrix}^{-1} \begin{bmatrix} \sum_{i=1}^n \int \rho_{1i}(\gamma) \hat{g}_i(\gamma) d\mathbb{Q}(\gamma) \\ \sum_{i=1}^n \int \rho_{2i}(\gamma) \hat{g}_i(\gamma) d\mathbb{Q}(\gamma) \end{bmatrix}.$$

We next verify the regularity conditions for the TSFLS estimator. First, note that both  $\rho_{1i}(\cdot)$  and  $\rho_{2i}(\cdot)$  are continuous on  $\Gamma = [0, 1]$  and  $\rho_i(\cdot)$  is linear with respect to  $(\theta_1, \theta_2)$ . Furthermore, if we suppose that  $\Theta$  is a compact and convex subset of  $\mathbb{R}^2$  containing  $\theta_* = (\pi_{\dagger}, \pi_{\dagger} \exp(\gamma_{\dagger}))$ , Assumption 2 holds trivially. Here, the uniqueness condition applies by virtue of the fact that  $\theta_{\dagger}$  is solely characterized by  $\pi_{\dagger}$  and  $\gamma_{\dagger}$ . Furthermore,  $A$  is given by

$$A = \begin{bmatrix} 1 + 2a\mu_2 + b\mu_4 & \mu_2 + 2a\mu_4 + b\mu_6 \\ \mu_2 + 2a\mu_4 + b\mu_6 & \mu_4 + 2a\mu_6 + b\mu_8 \end{bmatrix},$$

where  $\mu_n := \mathbb{E}[z_i^n]$ , and  $a$  and  $b$  are  $\int \exp(\gamma) d\mathbb{Q}(\gamma)$  and  $\int \exp(2\gamma) d\mathbb{Q}(\gamma)$ , respectively. Then  $A$  is positive definite by Cauchy-Schwarz and Assumption 4 holds. Second, with an IID data structure for  $(z_i, \varepsilon_i, \nu_i)'$ , the functional data set  $\{(\tilde{g}_i(\cdot), z_i)\}$  is IID also, and the measurability conditions in Assumption 6 hold by the continuity condition of  $\tilde{g}_i(\cdot)$  with respect to  $z_i, y_i, x_i, \gamma, \psi$ , and  $\delta$ . In particular, for each  $\gamma$ ,  $\tilde{g}_i(\gamma, \cdot)$  is differentiable with respect to  $\psi$  and  $\delta$ . Further,  $\Gamma = [0, 1]$  is a compact and convex set in the Euclidean real line, and we also suppose that the parameter spaces of  $\psi$  and  $\delta$  are compact and convex subsets of  $\mathbb{R}^2$  and  $\mathbb{R}$  containing  $\psi_{\dagger}$  and  $\delta_{\dagger}$ , respectively, so that Assumption 6 is satisfied. Third, we note that the upper bounds on the left sides of Assumptions 7(i-v) and 11(i-v) can be given in the form of  $(\zeta_0 + \zeta_1 z_i^2)(\kappa_0 + \kappa_1 y_i^2 + \kappa_2 z_i^2)$ . For example,  $|(\partial \rho_i(\gamma, \theta)) / (\partial \theta_1)| \leq (1 + \exp(1)z_i^2)$ , and we can let  $(\zeta_0, \zeta_1, \kappa_0, \kappa_1, \kappa_2)$  be  $(1, \exp(1), 1, 0, 0)$  for this case. Thus, if  $m_i$  is the maximum of the upper bounds, the  $L_2$  condition holds for  $m_i$  from the distributional condition for  $z_i, \varepsilon_i$ , and  $\nu_i$ , thereby verifying Assumptions 7 and 11. Fourth, note that the quasi-maximum likelihood estimator  $\hat{\xi}_n := (\hat{\psi}'_n, \hat{\delta}'_n)'$  is obtained by assuming a normal distribution for the error term, so that we can let  $s_i(\xi_*) = (x'_i u_i / \delta_*, u_i^2 / (2\delta_*^2) - 1 / (2\delta_*))'$ , and  $H = -\text{diag}[\mathbb{E}[x_i x'_i] / \delta_*, 1 / (2\delta_*^2)]$ , from which Assumption 8 holds, as White (1982) demonstrates. Furthermore, it is not difficult to provide estimators consistent for  $J$  and  $H$ . We can let  $\hat{J}_n$  and  $\hat{H}_n$  be their sample analogs, viz.,  $\hat{J}_n := n^{-1} \sum_{i=1}^n \hat{s}_i \hat{s}'_i$  and  $\hat{H}_n := -\text{diag}[n^{-1} \sum_{i=1}^n x_i x'_i / \hat{\delta}_n, 1 / (2\hat{\delta}_n^2)]$ , where  $\hat{s}_i := (x'_i \hat{u}_i / \hat{\delta}_n, \hat{u}_i^2 / (2\hat{\delta}_n^2) - 1 / (2\hat{\delta}_n^2))'$  and  $\hat{u}_i := y_i - x'_i \hat{\psi}_n$ . Then,  $\hat{J}_n$  and  $\hat{H}_n$  are consistent for  $J$  and  $H$  by the law of large numbers and the consistency of  $\hat{\xi}_n$  follows, verifying Assumption 10. Finally, we verify Assumption 9. As our main interest is in obtaining the null limit distribution of the tests defined below using the asymptotic distribution, we impose  $\pi_{\dagger} = 0$  and verify Assumption 9 to avoid the associated algebraic complexity. Some algebraic calculations show that

$$J = \begin{bmatrix} \frac{1}{\delta_*} & \frac{\mu_1}{\delta_*} & 0 \\ \frac{\mu_1}{\delta_*} & \frac{\mu_2}{\delta_*} & 0 \\ 0 & 0 & \frac{1}{2\delta_*^2} \end{bmatrix}, \quad K = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 + 2a\mu_2 + b\mu_4 & \mu_2 + 2a\mu_4 + b\mu_6 & 0 \end{bmatrix}, \quad \text{and}$$

$$B = 2\delta_*^2 \begin{bmatrix} 1 + 4a\mu_2 + 2(b + 2a^2)\mu_4 + 4ab\mu_6 + b^2\mu_8 & \mu_2 + 4a\mu_4 + 2(b + 2a^2)\mu_6 + 4ab\mu_8 + b^2\mu_{10} \\ \mu_2 + 4a\mu_4 + 2(b + 2a^2)\mu_6 + 4ab\mu_8 + b^2\mu_{10} & \mu_4 + 4a\mu_6 + 2(b + 2a^2)\mu_8 + 4ab\mu_{10} + b^2\mu_{12} \end{bmatrix}.$$

In Section 5, we conduct simulations by assuming  $\delta_{\dagger} = 1$  and a standard uniform distribution for  $\mathbb{Q}$ , from which it follows that  $a = \int_0^1 \exp(\gamma) d\gamma$ ,  $b = \int_0^1 \exp(2\gamma) d\gamma$ , and  $\delta_* = 1$  as  $\pi_{\dagger} = 0$ . Furthermore, from the DGP condition that  $z_i \sim \text{IID } U[0, 1]$ ,  $\mu_n = 1/(1 + n)$ . From these calculations the matrix  $C$  is positive definite. Finally, if  $\pi_{\dagger} = 0$ ,  $M = -K'$  and  $J = -H$  from the DGP condition, so that  $B_* = B - K'J^{-1}K$ , which is positive definite as  $C$  is positive definite. This verifies the regularity conditions for the TSFLS estimator.  $\square$

## 4 Inference on the Mean Function

Parametric specification of the conditional mean of functional data is particularly useful in inference. Instead of conducting inference over  $\Gamma$ , we can test a relevant hypothesis by estimating the unknown parameter  $\theta_*$  directly. In what follows we extend the standard analysis of Wald, Lagrange multiplier (LM), and quasi likelihood-ratio (QLR) test statistics to perform inference on the functional mean. Specifically, suppose the following hypotheses on the mean function are to be tested:

$$\mathbb{H}_o : h(\theta_*) = 0 \text{ versus } \mathbb{H}_a : h(\theta_*) \neq 0.$$

We assume that the function  $h(\cdot)$  is given and satisfies Assumption 12:

**Assumption 12. (Hypothesis):** (i)  $h : \Theta \mapsto \mathbb{R}^r$  is in  $\mathcal{C}^1(\Theta)$  with  $r \in \mathbb{N}$  and  $r \leq d$ ; and

(ii)  $D(\theta_*) := \nabla'_\theta h(\theta_*)$  has full rank  $r$ .  $\square$

Define the constrained FLS (CFLS) and constrained two-stage FLS (CTSFLS) estimators as

$$\check{\theta}_n^b := \arg \min_{\theta \in \Theta; h(\theta)=0} q_n(\theta) \quad \text{and} \quad \check{\theta}_n^\sharp := \arg \min_{\theta \in \Theta; h(\theta)=0} \hat{q}_n(\theta).$$

These restricted estimates are used in developing the test statistics. In view of the constrained optimization, the criteria  $q_n(\check{\theta}_n^b)$  and  $q_n(\check{\theta}_n^\sharp)$  cannot be smaller than  $q_n(\hat{\theta}_n)$  and  $\hat{q}_n(\tilde{\theta}_n)$ . In a similar way, we define  $\theta_{\dagger}$  to be the minimizer of  $q(\theta)$  under the same restriction, i.e.,  $\theta_{\dagger} := \arg \min_{\theta \in \Theta; h(\theta)=0} q(\theta)$ . The following Lemma establishes consistency of the CFLS and CTSFLS estimators and is repeatedly used in the limit theory of the test statistics.

**Lemma 1.** (i) Given Assumptions 1, 2, 3, and 12,  $\check{\theta}_n^b \rightarrow \theta_{\dagger}$  a.s.  $-\mathbb{P}$ , and  $\theta_{\dagger} = \theta_*$  under  $\mathbb{H}_o$ ; and

(ii) Given Assumptions 2, 6, 7, 8, 11, and 12,  $\check{\theta}_n^\sharp \rightarrow \theta_{\dagger}$  a.s.  $-\mathbb{P}$ , and  $\theta_{\dagger} = \theta_*$  under  $\mathbb{H}_o$ .  $\square$

### 4.1 Wald Test

We construct the usual Wald (1943) statistics as follows

$$\mathcal{W}_n^b := nh(\hat{\theta}_n)' \{ \hat{D}_n \hat{A}_n^{-1} \hat{B}_n \hat{A}_n^{-1} \hat{D}_n' \}^{-1} h(\hat{\theta}_n); \quad \text{and} \quad \mathcal{W}_n^\sharp := nh(\tilde{\theta}_n)' \{ \tilde{D}_n \tilde{A}_n^{-1} \tilde{B}_n \tilde{A}_n^{-1} \tilde{D}_n' \}^{-1} h(\tilde{\theta}_n)$$

where  $\hat{D}_n := D(\hat{\theta}_n)$ ,  $\tilde{D}_n := D(\tilde{\theta}_n)$  and all other notation is the same as in Section 3. The statistic  $\mathcal{W}_n^b$  is used for models without nuisance effects, whereas  $\mathcal{W}_n^\sharp$  is used for models with nuisance effects. The next result provides limit theory for these

Wald tests.

**Theorem 8.** (i) Given Assumptions 1, 2, 3, 4, 5, and 12,

(i.a)  $\mathcal{W}_n^b \stackrel{\Delta}{\sim} \mathcal{X}^2(r, 0)$  under  $\mathbb{H}_o$ , where  $\mathcal{X}^2(a, b)$  denotes a noncentral chi-square variable with degrees of freedom  $a$  and noncentrality parameter  $b$ ;

(i.b) for any sequence  $c_n \rightarrow \infty$  such that  $c_n = o(n)$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}[\mathcal{W}_n^b \geq c_n] = 1$  under  $\mathbb{H}_a$ ; and

(ii) Given Assumptions 2, 4, 6, 8, 9, 10, 11, and 12,

(ii.a)  $\mathcal{W}_n^\# \stackrel{\Delta}{\sim} \mathcal{X}^2(r, 0)$  under  $\mathbb{H}_o$ ; and

(ii.b) for any sequence  $c_n \rightarrow \infty$  such that  $c_n = o(n)$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{W}_n^\# \geq c_n) = 1$  under  $\mathbb{H}_a$ . □

The null limit distribution of the Wald test statistic is chi-squared with degrees of freedom  $r$ , where  $r$  is the rank of  $D_*$  as given in Assumption 12. The result follows easily from the fact that  $h(\hat{\theta}_n)$  and  $h(\tilde{\theta}_n)$  are asymptotically normal. On the other hand, when the null hypothesis is false, the Wald statistics diverge with probability one, giving consistency of the tests.

## 4.2 Lagrange Multiplier (LM) Test

The LM test statistics are defined as follows

$$\mathcal{LM}_n^b := \frac{n}{4} \nabla'_{\theta} q_n(\hat{\theta}_n^b) \hat{A}_n^{-1} \ddot{D}_n^{b'} \{ \ddot{D}_n^b \hat{A}_n^{-1} \hat{B}_n \hat{A}_n^{-1} \ddot{D}_n^{b'} \}^{-1} \ddot{D}_n^b \hat{A}_n^{-1} \nabla_{\theta} q_n(\hat{\theta}_n^b) \quad \text{and}$$

$$\mathcal{LM}_n^\# := \frac{n}{4} \nabla'_{\theta} \hat{q}_n(\hat{\theta}_n^\#) \tilde{A}_n^{-1} \ddot{D}_n^{\#'} \{ \ddot{D}_n^\# \tilde{A}_n^{-1} \tilde{B}_n \tilde{A}_n^{-1} \ddot{D}_n^{\#'} \}^{-1} \ddot{D}_n^\# \tilde{A}_n^{-1} \nabla_{\theta} \hat{q}_n(\hat{\theta}_n^\#),$$

where  $\ddot{D}_n^b := D(\hat{\theta}_n^b)$  and  $\ddot{D}_n^\# := D(\hat{\theta}_n^\#)$ . Here,  $\hat{B}_n$  and  $\tilde{B}_n$  can be replaced by other consistent estimators for  $B$  and  $B_*$ . For example, if we let

$$\ddot{B}_n^b := \frac{1}{n} \sum_{i=1}^n \int \int \nabla_{\theta} \rho_i(\gamma, \hat{\theta}_n^b) \{ g_i(\gamma) - \rho_i(\gamma, \hat{\theta}_n^b) \} \{ g_i(\tilde{\gamma}) - \rho_i(\tilde{\gamma}, \hat{\theta}_n^b) \} \nabla'_{\theta} \rho_i(\gamma, \hat{\theta}_n^b) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma}) \quad \text{and}$$

$$\ddot{B}_n^\# := \frac{1}{n} \sum_{i=1}^n \int \int \nabla_{\theta} \rho_i(\gamma, \hat{\theta}_n^\#) \{ g_i(\gamma, \hat{\xi}_n) - \rho_i(\gamma, \hat{\theta}_n^\#) \} \{ g_i(\tilde{\gamma}, \hat{\xi}_n) - \rho_i(\tilde{\gamma}, \hat{\theta}_n^\#) \} \nabla'_{\theta} \rho_i(\gamma, \hat{\theta}_n^\#) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma}),$$

it is clear that  $\ddot{B}_n^b$  and  $\ddot{B}_n^\#$  are both consistent for  $B$  under  $\mathbb{H}_o$ .

Asymptotic theory relies on the first-order derivatives of  $q_n$  and  $\hat{q}_n$  evaluated at the constrained estimates  $\hat{\theta}_n^b$  and  $\hat{\theta}_n^\#$ . Under regularity conditions we have

$$\nabla_{\theta} q_n(\hat{\theta}_n^b) = -\frac{2}{n} \sum_{i=1}^n \int \{ g_i(\gamma) - \rho_i(\gamma, \hat{\theta}_n^b) \} \nabla_{\theta} \rho_i(\gamma, \hat{\theta}_n^b) d\mathbb{Q}(\gamma) \quad \text{a.s.} - \mathbb{P} \quad \text{and}$$

$$\nabla_{\theta} \hat{q}_n(\hat{\theta}_n^\#) = -\frac{2}{n} \sum_{i=1}^n \int \{ g_i(\gamma, \hat{\xi}_n) - \rho_i(\gamma, \hat{\theta}_n^\#) \} \nabla_{\theta} \rho_i(\gamma, \hat{\theta}_n^\#) d\mathbb{Q}(\gamma) \quad \text{a.s.} - \mathbb{P}$$

as shown in Lemma 3(iii) in the Appendix. The following theorem then holds.

**Theorem 9.** (i) Given Assumptions 1, 2, 3, 4, 5, and 12, we have:

(i.a)  $\mathcal{LM}_n^b \stackrel{\Delta}{\sim} \mathcal{X}^2(r, 0)$  under  $\mathbb{H}_o$ ; and

(i.b) for any sequence  $c_n \rightarrow \infty$  such that  $c_n = o(n)$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{LM}_n^b \geq c_n) = 1$  under  $\mathbb{H}_a$ ;

(ii) Given Assumptions 2, 4, 6, 8, 9, 10, 11, and 12,

(ii.a)  $\mathcal{LM}_n^\# \stackrel{A}{\sim} \chi^2(r, 0)$  under  $\mathbb{H}_o$ ; and

(ii.b) for any sequence  $c_n \rightarrow \infty$  such that  $c_n = o(n)$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{LM}_n^\# \geq c_n) = 1$  under  $\mathbb{H}_a$ .  $\square$

Theorem 9 delivers the limit behavior of the LM statistics under the null and alternative hypotheses. The same null limit distributions apply as for the Wald statistic, because under  $\mathbb{H}_o$ , both  $\nabla_{\theta} q_n(\hat{\theta}_n^b)$  and  $\nabla_{\theta} q_n(\hat{\theta}_n^\#)$  are asymptotically normal with mean zero and covariance matrices that are consistently estimated by the weight matrices employed in construction of the LM statistics.

### 4.3 Quasi Likelihood Ratio (QLR) Test

The QLR statistics are defined for the models without and with nuisance effects as

$$\mathcal{QLR}_n^b := n\{q_n(\hat{\theta}_n^b) - q_n(\hat{\theta}_n)\} \quad \text{and} \quad \mathcal{QLR}_n^\# := n\{\hat{q}_n(\hat{\theta}_n^\#) - \hat{q}_n(\hat{\theta}_n)\}.$$

Approximating  $q_n(\cdot)$  (resp.  $\hat{q}_n(\cdot)$ ) via a second-order Taylor expansion yields the asymptotic distribution of  $\sqrt{n}(\hat{\theta}_n^b - \hat{\theta}_n)$  (resp.  $\sqrt{n}(\hat{\theta}_n^\# - \hat{\theta}_n)$ ), which is normal under  $\mathbb{H}_o$ . When  $\mathbb{H}_o$  is not true, this quantity is not bounded in probability, thereby distinguishing the null and alternative. But the QLR statistics do not have limiting chi-square distributions. Instead, their limit behavior under  $\mathbb{H}_o$  and  $\mathbb{H}_a$  is given in the following result.

**Theorem 10.** (i) Given Assumptions 1, 2, 3, 4, 5, and 12,

(i.a)  $\mathcal{QLR}_n^b \stackrel{A}{\sim} W' \{D_* A^{-1} D_*'\}^{-1} W$  under  $\mathbb{H}_o$ , where  $D_* := D(\theta_*)$ , and  $W \sim \mathcal{N}(0, D_* A^{-1} B A^{-1} D_*')$ ; and

(i.b) for any sequence  $c_n \rightarrow \infty$  such that  $c_n = o(n)$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{QLR}_n^b \geq c_n) = 1$  under  $\mathbb{H}_a$ ; and

(ii) Given Assumptions 2, 4, 6, 8, 9, 10, 11, and 12; and

(ii.a)  $\mathcal{QLR}_n^\# \stackrel{A}{\sim} W_*' \{D_* A^{-1} D_*'\}^{-1} W_*$  under  $\mathbb{H}_o$ , where  $W_* \sim \mathcal{N}(0, D_* A^{-1} B_* A^{-1} D_*')$ ; and

(ii.b) for any sequence  $c_n \rightarrow \infty$  such that  $c_n = o(n)$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{QLR}_n^\# \geq c_n) = 1$  under  $\mathbb{H}_a$ .  $\square$

The null limit distributions of the QLR test statistics differ from standard chi-squared theory because the asymptotic covariance matrices of the FLS and TSFLS estimators differ from the limits of the Hessian matrices of FSMSE and the information matrix equality fails.

The different null limit distributions of the QLR test statistics can be linked to the previous literature. Note that Vuong (1989) and White (1994) examine the null limit distribution of the likelihood-ratio test statistic by supposing that the maximum likelihood estimator has a sandwich-form asymptotic covariance matrix that typically follows from model misspecification and/or conditionally heteroskedastic errors, producing a noncentral chi-squared null limit distribution for their respective likelihood ratio test statistics. If the information matrix equality holds, so that if  $B = \sigma^2 A$ , the QLR test statistics can be converted to follow a chi-squared distribution under the null by dividing the QLR test statistics with a consistent estimator for  $\sigma^2$  under the null. This feature implies that the information matrix equality is necessary when constructing a likelihood ratio test for the limit theory to follow a chi-squared distribution under the null. Importantly,  $\hat{B}_n$  and  $\tilde{B}_n$  still play a critical role in applying the QLR test statistics. For even though computation of the QLR statistics does not rely on these matrices, they are needed to obtain critical values of the QLR tests.

## 5 Model Applications and Simulations

This section explains how the theory in Sections 3 and 4 for functional data relate to standard econometric analysis commonly used in applications. Four examples are given showing how the limit theory is applied in such settings. In the first application, we study model estimation conducted by the “smoothing first, then estimation” principle and compare parameter estimation obtained by FLS and random effect model estimation. In the second and third applications, we study Examples 1 and 2 in Monte Carlo experiments. The fourth application extends Example 1 to copula mixtures representing dependence structures between two variables. These applications are all developed within the functional data framework. Experiments are conducted to assess the adequacy of the limit theory in finite samples.

### 5.1 FLS and Random Effect Model Estimation

When a panel data set is available, the so-called “smoothing first, then estimation” principle can be applied to estimate the parameters associated with the population mean function. Specifically, if individual  $i$ 's state variable in period  $t$  is  $w_{it}$ , a panel dataset of the observations is the collection  $\{w_{it} : i = 1, 2, \dots, n; \text{ and } t = 1, 2, \dots, T\}$ . Such panel observations are discrete over time  $t$  for each  $i$ . It is often appealing to assume that the underlying individual  $i$ 's state path forms a set of discrete observations of a continuous curve over a certain time domain such as  $[0, T]$ . The corresponding interpretation of the discrete panel observations is that  $w_{it} = w_i(t)$  for the discrete time periods  $t = 1, 2, \dots, T$ , where the random function  $w_i(\cdot)$  denotes individual  $i$ 's full time path of states over the continuous time interval  $[0, T]$ . The full state path is unobserved. While a generating mechanism for the continuous process  $w_i(\cdot)$  might be constructed, it is often convenient to provide a direct empirical approximation by a sieve or local polynomial kernel method. For example, Zhang and Chen (2007) examined a smooth sample path approximated from  $\{w_{it} : t = 1, 2, \dots, T\}$  by local polynomial kernel estimation, and gave mild regularity conditions under which inference on the implied functional data  $\{g_i(\cdot) : i = 1, 2, \dots, n\}$  is asymptotically equivalent to that obtained from the discrete path  $\{w_{it} : i = 1, 2, \dots, n; t = 1, \dots, T\}$ .

Our model and econometric approach are motivated by the mean log income path application in Section 6. We therefore designed and conducted a simulation experiment to ascertain how the FLS approach can help to improve panel estimation. For this experiment we first generated the ‘unobserved’ functional data  $w_i(\cdot)$  according to the model

$$w_i(t) = \theta_{1*} + \theta_{2*}t + \theta_{3*}t^2 + \theta_{4*}t^3 + \theta_{5*}t^4 + \varepsilon_i(t), \quad (10)$$

where  $t \in [0, T]$  with  $T = 40$ , and  $\varepsilon_i(\cdot)$  is an Ornstein-Uhlenbeck (OU) error process generated by the stochastic differential equation

$$d\varepsilon_i(t) = -\kappa_{\dagger}\varepsilon_i(t)dt + \sigma_{\dagger}d\mathcal{W}_i(t),$$

where  $\mathcal{W}_i(\cdot)$  is a collection of standard Wiener process that are identically and independently distributed across individuals. For the simulation, we set  $\kappa_{\dagger} = \sigma_{\dagger} = 1$ . The panel observations are obtained as  $\{w_{it}\}$  by evaluating  $w_i(t)$  at the discrete time periods  $t = 1, 2, \dots, 40$  for each individual  $i = 1, 2, \dots, n$ , and the functional observations  $\{g_i(\cdot) : i = 1, 2, \dots, n\}$  for the FLS approach are obtained by smoothing the data  $\{w_{it}\}$  using a local polynomial kernel. After generating the data, the random effects model is estimated according to Wooldridge (2010) by using the intercept and polynomial time trends as regressors, while the FLS estimator is estimated by letting the adjunct probability measure  $\mathbb{Q}$  be a uniform distribution on  $[0, 40]$ . The assumption of the uniform distribution is made to correspond to the supposition that the researcher has no particular prior information concerning

how the null hypothesis of interest is violated when the FLS estimator is used to test a hypothesis.

Some further remarks on this experiment are warranted. First, our empirical application is concerned with income processes, which may be well suited to local polynomial estimation. According to Zhang and Chen (2007), local polynomial estimation provides better estimates of functional observations from discrete observations if the hidden functional observations are smooth. Such latent functions are amenable to the use of their optimal bandwidth if  $T$  and  $n$  are large and the latent function is smooth then a polynomial approximation can be expected to work especially well, as polynomials are dense in the space of continuous functions on a compact set and typically provide good smooth approximants to underlying smooth functions. In addition, the literature on income processes provides ample evidence that income processes are smooth series. Lillard and Weiss (1979) hypothesize that individual labor income processes exhibit deterministically growing individual-specific income growth with a stationary component. Guvenen (2007) calls this property the heterogeneous income profile (HIP) hypothesis contrasting it to the so-called restricted income profile (RIP) that supposes labor income profiles are homogeneous among individuals but subject to persistent aggregate income shocks. Guvenen (2007) conducts tests to validate each hypothesis and concludes that several notable features of consumption are more compatible with HIP than RIP; and Guvenen (2009) provides further evidence for HIP. As another example, Lintner (1956) conducted interviews with managers from twenty-eight companies that revealed a potential smoothing characteristic in the behavior of dividend income. From these interviews Lintner found that firms are reluctant to announce dividend changes which they might be obliged to reverse in the future, leading him to argue that dividends are adjusted in response to non-transitory earnings changes with the goal of achieving a long-run target payout ratio. Other empirical work supports this view (e.g. Fama and Blahnik, 1968; Marsh and Merton, 1987; Garrett and Priestley, 2000; Brav, Graham, Harvey, and Michaely, 2005; and Andres, Betzer, Goergen, and Renneboog, 2009). All these studies provide some support for the use local polynomial estimation as a suitable methodology to estimate income processes.

Second, the OU process with  $\kappa_* > 0$  is asymptotically stationary with zero mean and nondifferentiable everywhere almost surely so that the function  $w_i(\cdot)$  is the composition of a smooth polynomial trend and a nondifferentiable stochastic process. Therefore, smooth local polynomial kernel estimation is inevitably biased for  $w_i(\cdot)$ . Nonetheless, the procedure produces functional data observations that capture the polynomial drift and enable estimation of the unknown parameters by FLS. The numerical performance of this procedure can be compared with that of standard random effect model estimation. The experiment will help to demonstrate how the FLS estimation methodology has the capacity dominate random effect model estimation even when the functional data are constructed without retrieving the exact continuous path features of the underlying data generating process. Third, the simulation experiment was designed to follow the empirical analysis in Section 6 in which the conditional population mean of the continuous income path is estimated from 40 year annual incomes of white male and female workers in the U.S. The simulated panel data is therefore structured in a similar way to the actual panel data used in Section 6 and the model estimated in the simulation exercise matches that of the empirical section, which is specified to examine the income path hypothesis posited in the early research by Mincer (1958, 1974). A further motivation for (10) stems from the fact that polynomials are dense in the space of continuous functions on a compact interval and the Stone-Weierstrass theorem demonstrates that any such continuous function can be uniformly approximated by a polynomial function of high enough degree. Past empirical research suggests that the fourth degree polynomial function in (10) is sufficient for capturing the main characteristic behavior of labor income processes (e.g., Murphy and Welch, 1990; Cho and Phillips, 2018a).

<<<<<< Insert Table 5.1. >>>>>>

Table 5.1 reports the bias, root mean square error (RMSE), and mean absolute percentage error (MAPE) of the estimators in



the random effects model and the FLS approach, computed from  $R = 5,000$  replications. For each  $j = 1, 2, 3$ , these summary statistics are computed by

$$\text{bias}(\hat{\theta}_{jn}) = \frac{1}{R} \sum_{i=1}^R (\hat{\theta}_{jn}^{(i)} - \theta_{j*}), \quad \text{RMSE}(\hat{\theta}_{jn}) = \sqrt{\frac{1}{R} \sum_{i=1}^R (\hat{\theta}_{jn}^{(i)} - \theta_{j*})^2}, \quad \text{and} \quad \text{MAPE}(\hat{\theta}_{jn}) = \frac{1}{R} \sum_{i=1}^R \left| 1 - \frac{\hat{\theta}_{jn}^{(i)}}{\theta_{j*}} \right| \times 100,$$

where  $\hat{\theta}_{jn}^{(i)}$  is the estimate in the  $i$ -th replication. The parameter values  $\theta_{j*}$  in (10) are chosen as the empirical parameter estimates obtained in Section 6 by FLS for the male sample group with a Bachelor degree. For the mean function, we explore several different specifications, covering the quadratic, cubic, quartic, and restricted quartic models employed in Section 6. Accordingly, the data are generated by imposing some restrictions on the parameter vector  $\theta_*$  depending on the specification. For example, the data for the quadratic model are generated by setting  $\theta_{4*}$  and  $\theta_{5*}$  to zero, and the parameters are estimated based on the quadratic specification of the mean function. For simplicity, the parameter  $\alpha$  is set to be  $1/60$  in the restricted quartic model as in Murphy and Welch (1990). We do not report the bias, RMSE, and MAPE of  $\hat{\theta}_{4n}$  and  $(\hat{\theta}_{4n}, \hat{\theta}_{5n})$  for the cubic and quartic models, respectively in Table 5.1 as they all have values close to zero.

As shown in Table 5.1, the FLS approach generally provides better performance than random effect model estimation. The bias, RMSE and MAPE of the parameter estimates are all smaller in the FLS approach in the simulation models examined. The improvement is more marked in the parameter estimates of the coefficients for lower degree polynomial terms and when the sample size is small. In addition, we have found in experiments not reported here that these simulation results are robust to different simulation environments. In particular, when the OU process is replaced by the Gaussian processes obtained while testing for Hermitian, exponential, Weibull mixtures (e.g., Cho and White, 2007 and 2010) or Cox-Ingersoll-Ross process, we continue to find that the FLS estimator outperforms the random effects estimation. These findings indicate the gains that can be achieved in use of the FLS estimation over random effect model estimation.

In addition to these simulations, some robustness exercises were conducted. All the models in this section were also fitted using a general beta distribution for the adjunct probability measure under the same model assumptions as those given in Table 5.1. These simulations produced broadly similar findings that did not modify any of the qualitative results obtained with a uniform measure. The performance of the FLS estimator with a fitted beta distribution for the adjunct measure did not therefore show any noticeable gains over FLS, at least for the models considered.

## 5.2 Example 1: Distribution Specification Tests – Continued

We next considered the performance of our Wald, LM, and QLR test statistics. Some of their key features are now summarized in relation to the test on the supremum over  $\gamma$  statistic,  $C(\alpha)$ . Critical values of the  $C(\alpha)$  test are obtained by explicitly exploiting the functional form of  $g_i(\cdot)$  in the construction of that statistic, making its application less convenient in practice. Making the correct model assumption assumes correct knowledge of the functional form of  $g_i(\cdot)$ , which raises specificity and reduces the range of applicability of the maximal test statistic. The null limit distribution of  $C(\alpha)$  test statistic is typically obtained by simulation or bootstrapping that exploits the covariance kernel structure of a Gaussian stochastic process derived by assuming that the mixture assumption is correct. By contrast, our test statistics do not exploit this feature and our test statistics allow for model misspecification for  $\mu$ . In consequence, our statistics lead to tests whose power may not be as great as that of the  $C(\alpha)$  statistic under the correct mixture model assumption; but otherwise there is no clear power ordering and our test statistics may provide more powerful diagnostics as well as computational ease in various circumstances such as misspecification.

For the simulation exercise, we specify the same  $\rho(\cdot)$  as before and generate  $x_i$  according to the mixture assumption, so that we can fix the experimental framework and compare the  $C(\alpha)$  test with the three other tests. We proceed according to the following plan. First, we test

$$\mathbb{H}_o : (\theta_{1*}, \theta_{2*}) = (0, 0) \quad \text{versus} \quad \mathbb{H}_a : (\theta_{1*}, \theta_{2*}) \neq (0, 0)$$

using the three test statistics  $\mathcal{Q}_n^b$ ,  $\mathcal{LM}_n^b$ , and  $\mathcal{QLR}_n^b$ . For this hypothesis, Assumption 12 is trivially satisfied. Note that these applications do not require that the mixture assumption is correct for  $x_i$ . That is,  $\widehat{B}_n$  is estimated irrespective of whether  $\mathcal{M}$  is correct for  $\mu(\cdot)$  or not. In addition, we also test  $\theta_{1*}$  along with  $\theta_{2*}$  as it provides useful information for the DGP. If  $\theta_{1*} \neq 0$ , it is a strong signal for that the given mixture assumption is not valid for  $x_i$ . Second, we test  $\pi_{\dagger} = 0$  by the  $C(\alpha)$  test statistic. Note that under this DGP condition, if  $\pi_{\dagger} = 0$ ,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n g_i(\cdot) \Rightarrow \mathcal{G}(\cdot)$$

such that for each  $\gamma$  and  $\tilde{\gamma}$ ,  $\mathbb{E}[\mathcal{G}(\gamma)\mathcal{G}(\tilde{\gamma})] = \mathbb{E}[g_i(\gamma)g_i(\tilde{\gamma})]$ . Therefore, the asymptotic critical values of the  $C(\alpha)$  test are given by the distribution of  $\sup_{\gamma \in \Gamma} \mathcal{G}(\gamma)$ , and Cho and White (2010) show that the distribution of  $\mathcal{G}(\cdot)$  is identical to  $\tilde{\mathcal{G}}(\cdot)$ , where for each  $\gamma$ ,

$$\tilde{\mathcal{G}}(\gamma) := \sum_{k=2}^{\infty} \left[ \frac{(\gamma-1)^4}{\gamma^2(2\gamma-1)} \right]^{-1/2} \left( \frac{\gamma-1}{\gamma} \right)^k Z_k,$$

and  $Z_k \sim \text{IID } N(0, 1)$ , so that the asymptotic critical values can be delivered by repeatedly simulating  $\sup_{\gamma \in \Gamma} \tilde{\mathcal{G}}(\gamma)$ .

Table 5.2 displays the size and power of the Wald, LM and QLR test statistics studied in Section 4 along with the  $C(\alpha)$  test. Throughout the experiment, the FLS estimator is estimated by Gauss-Legendre numerical quadrature.  $\Gamma$  is chosen to be the interval  $[1.5, 2.5]$ , and as long as  $\underline{\gamma} > 1$ , this interval is arbitrarily selected to accommodate the fact that the researcher may not have information on the underlying DGP. We also let the adjunct probability measure be the uniform probability measure on  $\Gamma$  and consider sample sizes  $n = 25, 50, 100, 300$ , and 500. The nominal levels are fixed at 1%, 5%, and 10%. In the level panel of Table 5.2, we observe that the rejection rates of the three test statistics all approach nominal levels as the sample size increases. Under the alternative with  $\gamma_{\dagger} = 2$ , power is computed through 5,000 replications with the same sample sizes, but with the nominal level fixed at 0.05. In particular, we examine power of the tests by letting  $\pi_{\dagger}$  vary over the range  $\{0.1, 0.2, 0.3, 0.4, 0.5\}$ . The rejection rates tend to be larger as we move  $\pi_{\dagger}$  further from zero; and when the sample size increases, rejection rates approach unity for fixed  $\pi_{\dagger}$ . In addition, the power of the  $C(\alpha)$  test is overall greater than that of each of the three test statistics, as expected from the correct mixture assumption. Despite the greater power of the  $C(\alpha)$  test under these conditions, we emphasize that its implementation requires considerably more effort, a correct model assumption, and simulated asymptotic critical values, which rely on simulating the correct covariance kernel structure, or use of a bootstrap approach. Tests based on the FLS estimator are by comparison quite straightforward to implement.

<<<<<< Insert Table 5.2. >>>>>>

### 5.3 Example 2: Inference on Random Coefficients – Continued

We next conduct simulations using the functional data affected by the nuisance parameter estimation error:  $\widehat{g}_i(\gamma) := \tilde{g}_i(\gamma, \widehat{\delta}_n, \widehat{\psi}_n)$ , and we test the following hypotheses:

$$\mathbb{H}_o : \theta_* := (\theta_{0*}, \theta_{1*})' = 0 \quad \text{versus} \quad \mathbb{H}_a : \theta_* \neq 0.$$

Here, we note that

$$\theta_{0*} = -\pi_{\dagger} \exp(\gamma_{\dagger}) \mathbb{E}[z_i^2] \quad \text{and} \quad \theta_{1*} := \pi_{\dagger} \exp(\gamma_{\dagger}),$$

so that if  $\pi_{\dagger} = 0$ ,  $\mathbb{H}_o$  holds, whereas  $\pi_{\dagger} \neq 0$  under  $\mathbb{H}_a$ .

In this FDA framework, we conduct simulations by applying the theorems for the models with nuisance effects in Section 3. We let  $(\psi_{1\dagger}, \psi_{2\dagger}, \gamma_{\dagger}, \delta_{\dagger}) = (1, 1, 0.5, 1)$  under both the null and alternative hypotheses, and set  $\pi_{\dagger} = 0$  under the null. For the alternative DGP, we consider various values for  $\pi_{\dagger}$ , viz., 0.01, 0.02, 0.03, 0.04, and 0.05. Next, let  $z_i \sim \text{IID } U[0, 1]$  and  $(\nu'_i, \varepsilon_i)' \sim \text{IID } \mathcal{N}(0, I_3)$ , so that the model for  $\mu(\cdot)$  is correctly specified. In addition, we let the adjunct probability measure be uniformly distributed over  $[0, 1]$  with  $\Gamma = [0, 1]$ . For the Wald and QLR test statistics, we use  $\tilde{A}_n$  and  $\tilde{B}_n$  when computing the test statistics, whereas  $\tilde{A}_n$  and  $\tilde{B}_n^{\#}$  are used in calculating the LM test statistic.

<<<<<< Insert Table 5.3. >>>>>>

Table 5.3 displays the empirical rejection ratios obtained for the Wald, LM and QLR statistics in testing the reformulated model and hypotheses. The results show that the null rejection rates are close to nominal levels for all three test statistics when the sample size is large. The Wald and LM test statistics tend to be slightly oversized when the sample sizes are small, whereas the QLR test statistic perform better in size control. For power analysis, we consider models with  $\pi_{\dagger} = 0.01, 0.02, 0.03, 0.04$ , and 0.05, fixing the nominal significance level at 5%. Evidently, the rejection rates turn out to be dependent on sample size for each value of  $\pi_{\dagger}$  in all three test statistics. As expected, the empirical rejection rates increase as  $\pi_{\dagger}$  or  $n$  increase. Overall, the QLR test shows better performance than the Wald and LM tests.

For comparison we display test results from Breusch and Pagan's (1979) conditional heteroskedasticity test statistic, denoted  $\mathcal{BP}_n$  in Table 5.3. If the sample size is small,  $\mathcal{BP}_n$  controls size better than the Wald and LM test statistics, but the QLR test statistic outperforms  $\mathcal{BP}_n$ . Under the alternative, the powers of the Wald, LM, and QLR tests are all superior to  $\mathcal{BP}_n$ .

## 5.4 Inference on the Homogeneity of Dependence Structure

As a further application we develop tests for heterogeneity in dependence structures by applying the mixture model assumption. For this purpose, suppose a researcher observes IID observations of multiple variables. Empirical interest often lies in determining whether the dependence structure among these variables is homogeneous. Even though the univariate marginals for each variate may remain constant over the whole population, observations can still be heterogeneous due to different dependence structures.

In what follows we provide test statistics to detect violations of homogeneity using finite copula mixtures. The Sklar theorem (1959) is useful for this purpose as it conveniently separates information on the univariate marginals from the joint distribution by means of the copula function. There is now a vast literature demonstrating the use of copulas for studying dependence structures (e.g., Nelsen, 2007; Joe, 2014; and the references therein) and many studies applying mixture copula models (e.g., Dias and Embrechts, 2004; Chen and Fan, 2006; Hu, 2006; Lai, Chen, and Gerlach 2009; Diks, Panchenko, and van Dijk 2010; Zimmer, 2012; Kosmidis and Karlis, 2016; Loaiza-Maya, Smith, and Maneesoonthorn 2018). But methods of inference concerning homogeneity in dependence structures based on finite mixture copula models has, to the best of our knowledge, so far not been addressed.

We proceed by considering a mixture of two distinct bivariate copula component densities  $c_1$  and  $c_2$  with parameter vectors  $\gamma_{1\dagger}$  and  $\gamma_{2\dagger}$ , respectively. That is, for  $(u, v) \in [0, 1]^2$ ,

$$c(u, v; \pi_{\dagger}, \gamma_{1\dagger}, \gamma_{2\dagger}) = (1 - \pi_{\dagger})c_1(u, v; \gamma_{1\dagger}) + \pi_{\dagger}c_2(u, v; \gamma_{2\dagger})$$

with  $\pi_{\dagger} \in [0, 1]$ . More generally, each component density can be of any dimension, and a mixture with more component densities can be considered. For brevity we focus on the simple prototypical model above. We suppose that IID pairs  $\{(x_i, y_i)\}_{i=1}^n$  have marginal distributions given as  $F_X$  and  $F_Y$ , respectively. Inference concerning the homogeneity of the dependence structure can naturally be conducted by examining the null hypothesis that  $\pi_{\dagger} = 0$  (or  $\pi_{\dagger} = 1$ ) with function  $g_i$  given by

$$g_i(U_i, V_i; \gamma_1, \gamma_2) := \frac{c_2(U_i, V_i; \gamma_2) - c_1(U_i, V_i; \gamma_1)}{c_1(U_i, V_i; \gamma_1) \sqrt{c^*(U_i, V_i; \gamma_1, \gamma_2) - 1}},$$

where  $U_i := F_X(x_i)$ ,  $V_i := F_Y(y_i)$ , and

$$c^*(u, v; \gamma_1, \gamma_2) := \int_0^1 \int_0^1 \frac{c_2^2(u, v; \gamma_2)}{c_1(u, v; \gamma_1)} dudv.$$

Here,  $g_i$  is derived by applying the  $C(\alpha)$  test principle for testing  $\pi_{\dagger} = 0$ . The derivation of  $g_i$  is reported in the Appendix.

A practical challenge arises from the fact that the univariate marginals  $F_X$  and  $F_Y$  are typically unknown to researchers. In the first stage estimation, we, therefore, approximate  $U_i$  and  $V_i$  by  $\hat{U}_i = \hat{F}_X(x_i)$  and  $\hat{V}_i = \hat{F}_Y(y_i)$ , respectively, using estimates of the marginal distributions in a fashion similar to the inference function for marginals (IFM) approach (Joe and Xu, 1996; Joe, 2001). To apply the results in Section 4.2, we construct the functional data as follows

$$\hat{g}_i(\gamma_1, \gamma_2) := \frac{c_2(\hat{U}_i, \hat{V}_i; \gamma_2) - c_1(\hat{U}_i, \hat{V}_i; \gamma_1)}{c_1(\hat{U}_i, \hat{V}_i; \gamma_1) \sqrt{c^*(\hat{U}_i, \hat{V}_i; \gamma_1, \gamma_2) - 1}}.$$

A leading example is the case when  $c_1(\cdot)$  is the density implied by the independence copula. If so, the functional form of  $g_i(\cdot)$  is further simplified as

$$\hat{g}_i(\gamma_2) := \frac{c_2(\hat{U}_i, \hat{V}_i; \gamma_2) - 1}{\sqrt{c^*(\hat{U}_i, \hat{V}_i; \gamma_2) - 1}}, \quad \text{where} \quad c^*(u, v; \gamma_2) = \int_0^1 \int_0^1 c_2^2(u, v; \gamma_2) dudv.$$

For our simulation experiments we use the Farlie-Gumbel-Morgenstern (FGM) copula for  $c_2(\cdot)$ , which enables a closed form solution for the relevant integrals. The population mean function is then straightforwardly derived as  $\mu(\pi_{\dagger}, \gamma_{2\dagger}) := \frac{1}{3}\pi_{\dagger}\gamma_{2\dagger}$ , leading to a simple linear model for the mean function, viz.,

$$\rho(\gamma, \theta_1, \theta_2) := \theta_1 + \theta_2\gamma,$$

where  $\gamma \in \Gamma := [0, 1]$ . Note that if  $\pi_{\dagger} = 0$ ,  $\mu(\pi_{\dagger}, \cdot) \equiv 0$ , and  $\rho(\cdot, \theta_{1*}, \theta_{2*}) \equiv 0$  if and only if  $(\theta_{1*}, \theta_{2*})' = (0, 0)'$ . We, therefore, specify the null and alternative hypotheses as follows:

$$\mathbb{H}_0 : (\theta_{1*}, \theta_{2*}) = (0, 0) \quad \text{vs.} \quad \mathbb{H}_a : (\theta_{1*}, \theta_{2*}) \neq (0, 0).$$

Although the intercept term is known to be zero, it is estimated and simulations are conducted based upon the multiple parameter estimators, enabling us to define test statistics with multiple parameters in the simulations.

The following computational algorithm is used in the simulations. First, we generate random samples using the FGM copula with marginals  $\mathcal{N}(0, 1)$  and  $\mathcal{N}(0, 5)$  for  $x_i$  and  $y_i$ , respectively, and we estimate the means and variances of  $x_i$  and  $y_i$  by maximum likelihood to obtain  $\hat{U}_i := \Phi(x_i, \hat{\mu}_{x,n}, \hat{\sigma}_{x,n}^2)$  and  $\hat{V}_i := \Phi(y_i, \hat{\mu}_{y,n}, \hat{\sigma}_{y,n}^2)$ , where  $\Phi(\cdot, \mu, \sigma^2)$  signifies the normal

distribution function with mean  $\mu$  and variance  $\sigma^2$ , and  $(\hat{\mu}_{x,n}, \hat{\sigma}_{x,n}^2)$  and  $(\hat{\mu}_{y,n}, \hat{\sigma}_{y,n}^2)$  are the corresponding maximum likelihood estimates obtained from  $x$  and  $y$  samples. Second, we fix the copula parameter  $\gamma_{2\ddagger}$  at 0.9, and the adjunct probability measure is assumed to be uniformly distributed on  $\Gamma := [0, 1]$ . With this framework, we conduct independent experiments with 5,000 replications using data samples with  $n = 25, 50, 100, 300, 500$ , and 1,000.

<<<<<< Insert Table 5.4. >>>>>>

Table 5.4 reports the empirical rejection rates, giving size and power of our tests for dependence. As before, the rejection rates given in the size panel are computed with fixed nominal levels of 1%, 5%, and 10%. Table 5.4 suggests that when the null hypothesis is true, the rejection rates approach nominal levels as the sample size increases. More specifically, the rejection rates of the Wald and LM test statistics are close to the nominal levels even with sample sizes as small as  $n = 50$  or 100, whereas those of the QLR test are somewhat oversized at these nominal levels and sample sizes. For power analysis we let  $\pi_{\ddagger}$  be 0.1, 0.2, 0.3, 0.4, and 0.5 and fix the nominal level of significance to 5%. A clear tendency for the rejection rates to rise with  $\pi_{\ddagger}$  is apparent; and when  $\pi_{\ddagger}$  is constant, the rejection rates of the three test statistics all approach unity as the sample size increases.

For comparison we also display the test results obtained from the Cramér–von Mises test and the Kolmogorov–Smirnov test for an independence copula, in which the statistics are computed using the  $l^2$  and  $l^\infty$  distances between the independence copula  $c_0$  and the empirical copula  $c_n$ , viz.,  $\sqrt{n} \|c_n - c_0\|_p$  with  $p = 2$  and  $\infty$ . In the table, these statistics are denoted  $\mathcal{CM}_n$  and  $\mathcal{KS}_n$ . As the null limit distributions of these test statistics are model dependent, the critical values are obtained by the bootstrap based on sampling with replacement (e.g., Efron, 1979; Fermanian, 2004). Table 5.4 shows that the Wald, LM and QLR tests generally perform better than the Cramér–von Mises and Kolmogorov–Smirnov tests.

## 6 Empirical Analysis on the Income Path

Lifetime earning trajectories have attracted great interest among labor economists, leading to the early earnings function pioneered by Mincer (1958, 1974), which revealed that earnings typically rise at a diminishing rate over a lifetime, justifying the use a quadratic form over work experience years in regression specifications. Since then, quadratic specifications with respect to years of work experience have been a popular component of empirical wage equations in the literature (e.g., Bhuller, 2017; Barth, Davis, and Freeman, 2018; Magnac, Pistoletti, and Roux, 2018 for recent studies). On the other hand, Murphy and Welch (1990) explored cubic and quartic specifications for wage equations, showing that a restricted quartic specification provided smaller MSE compared to quadratic specifications. Later, Katz and Murphy (1992), Autor, Katz, and Krueger (1998), and Lemieux (2006) adopted quartic specifications in their empirical work. Cho and Phillips (2018b) examined functional form specifications of the wage equation with respect to work experience years using sequential testing and found that functional form can be sensitive to the presence of other explanatory variables in the regression. In another study, Heckman, Lochner, and Todd (2006) showed that the quadratic model is empirically misspecified for recent wage data, but that relaxation of the quadratic model to a quartic specification does not dramatically change the empirical economic implications of the quadratic model.

This section applies the techniques of the current work to study the mean log income path (MLIP) as a function of work experience. Using several different formulations for the MLIP, we examine whether the overall shape of the MLIP differs in significant ways between genders and amongst education levels. In cases where the difference is significant, we further study how gender and education affect the MLIP, providing some empirical insights to enhance understanding of how different income profiles arise according to gender and education.

Our analysis differs from existing work based upon the estimation of Mincer (1958, 1974) wage equations in two ways. Primarily, our approach uses functional observations that have never before, to our knowledge, been considered in this literature. Next, although we estimate the MLIP to identify gender and education effects in parallel to methodological developments in the wage equation literature, our empirical results here are obtained by recasting the methodology into a form suited to functional data to single out how gender and education influence the overall shape of the MLIP. This functional approach has several advantages over the existing methodology. First, as shown in the simulation in Sections 5.1, the FLS using functional observations outperforms the usual panel approach. Second, the FDA transforms what would be the analysis of panel data with complicated temporal dependence structures into temporal curve analysis using IID observations. Accordingly, simple limit theories for an iid framework can be adopted. A final advantage is that the functional data approach can be applied even when the data are observed at random rather than at equispaced discrete times in the continuous domain. In contrast, econometric analysis using unbalanced panel observations is complicated and can be conducted by imposing assumptions on the generating mechanism that may not be justified for the particular data set in hand.

As it turns out in our empirical application, curve shape is not affected by gender and educational differences if each individual's income path is properly scaled by the individual's integrated log income path (LIP) over the work experience years, implying that different gender and education levels lead to different income paths for individuals, but individual LIPs are proportional to each other between genders and education levels. In short, appropriately scaled LIPs are shown to be unaffected by different gender and education levels.

For the empirical analysis we analyze income data obtained from the Continuous Work History Sample (CWHS) database that provides the income variable as annual labor income before taxes. The data include 39 years of income records of full time white male and female workers in the U.S. who were born between 1960 and 1962 and had tax records for at least 39 years. Based on these observations we construct each individual's LIP using the local polynomial kernel (Zhang and Chen, 2007) defined over work experience years from 0 to 40 years and then subdivide the entire sample into different groups based on individual gender and education level. For the latter the subdivisions are no college education, Bachelor's degree, Master's degree, and Doctoral degree. For the education level groups (in the order given above), we have 673, 2,828, 539, and 323 income path samples for males and 837, 1,624, 469, and 418 income path samples for females.

We conduct our empirical analysis by controlling for worker job mobility. Earlier literature has noted that Mincer's (1958, 1974) quadratic equation gives a good local approximation of the wage equation with respect to the work experience years. For example, Light and Ureta (1995) observe that the empirical wage profiles are more heterogenous in early career experience than those of workers reaching ages in the forties and fifties, in addition to their more rapidly increasing wage levels during this period of the life cycle. This phenomenon in the data is explained through more frequent job interruptions over the early career experience years, particularly for female workers. The finding is used to link gender wage gaps to the work experience variable. In addition to the work of Light and Ureta (1995), many other studies point out that wage profiles during early career years differ considerably from those of generations with longer work experience. Geographic changes and job mobility dominate among the young generation and young workers may have different perspectives on lifetime wage profiles from older workers when firms bargain over wage-employment packages with labor unions (e.g. Mincer and Jovanovic, 1981; Huizinga, 1990). These characteristics of the labor market motivate the treatment of individual income profiles as a composite of many different income profiles.

We, therefore, separately examine the income profiles defined over the entire working lifetime to profiles based on more mature work experience years. Specifically, we proceed by first focusing on the income paths over the full 0 – 40 work experience

years and examine how different gender and education levels affect the income paths. Next, we consider income paths over the 10 – 40 years work experience years and examine gender and education effects on these paths. In doing so, the first 10 years of the income paths are removed from the original paths to accommodate the arguments made by Mincer and Jovanovic (1981). These authors found empirical evidence from the NLS and MID panel data that differences in job mobility during the first 10 years of work experience do not predict long-run differences in earnings, implying that income paths during the first 10 years are likely generated by a mechanism that differs from that determining income paths in mature career experience years. As will become clear, we use this empirical separation as a platform to highlight different gender and education effects on income paths.

## 6.1 Inference on the Mean of Log Income Path on the Whole Work Experience Years

This section reports estimates of the MLIP over the entire working lifetime based on quadratic, cubic, and quartic models. We also estimate the restricted model posited by Murphy and Welch (1990), compare the results, and discuss inferences on the MLIP that are implied by these estimates.

<<<<<< Insert Figure 1. >>>>>>

Figure 1 shows the estimated shapes of the MLIPs with respect to work experience over 0 – 40 years for groups classified according to gender and education levels, which are implied by the quadratic, cubic, and quartic models. The red lines in the Figure denote pointwise MLIPs and the dashed lines present the 80% bootstrap confidence bands of the pointwise MLIPs that were obtained by resampling the functional observations with replacement. Along with the pointwise MLIPs are shown the MLIP curves corresponding to the quadratic, cubic, and quartic models. The quartic model, for example, is specified for each group as

$$\rho(\gamma, \theta_1, \theta_2, \theta_3, \theta_4, \theta_5) = \theta_1 + \theta_2\gamma + \theta_3\gamma^2 + \theta_4\gamma^3 + \theta_5\gamma^4.$$

After estimating the unknown parameters by FLS the fitted curves  $\rho(\cdot, \hat{\theta}_{1n}, \hat{\theta}_{2n}, \hat{\theta}_{3n}, \hat{\theta}_{4n}, \hat{\theta}_{5n})$  are shown for the entire working lifetime in the Figure. Here, the unknown parameters are estimated by letting the adjunct probability measure  $\mathbb{Q}$  be the uniform distribution on  $\Gamma$ , implying that

$$(\hat{\theta}_{1n}, \hat{\theta}_{2n}, \hat{\theta}_{3n}, \hat{\theta}_{4n}, \hat{\theta}_{5n})' = \left[ \int_{\Gamma} \gamma \gamma' d\gamma \right]^{-1} \left[ n^{-1} \sum_{i=1}^n \int_{\Gamma} \gamma g_i(\gamma) d\gamma \right],$$

where  $\gamma := (1, \gamma, \gamma^2, \gamma^3, \gamma^4)'$ . Note that the FLS estimator is well defined from the fact that the inverse matrix is nonsingular, implying that the model is identified. This representation is justified on the basis that there is no particular prior information regarding which  $\gamma$  violates the null hypotheses as given below, and the earlier simulation evidence in Section 5.1 does not lead to a substantially more efficient FLS estimator by generalizing the adjunct probability to the best distribution. We thus assume an equal chance for the null to be violated for each  $\gamma$  so that alternative behavior of the tests is reflected by the uniform distribution.

Although all estimated MLIPs in Figure 1 lie inside the 80% confidence bands of the pointwise MLIPs, the fitted MLIP curves are evidently different from the pointwise MLIPs, implying that it is difficult to estimate the MLIPs by quadratic, cubic, and quartic model specification uniformly over the full range of years of work experience. In particular, the pointwise MLIPs over the first 10 work experience years are clearly different from those implied by the model estimates. This may indicate, for example, that high job mobility during the first decade of a working lifetime produces income profiles different from those over the remaining years of work experience, corroborating the argument of Mincer and Jovanovic (1981). In addition to the MLIPs

implied by the model estimates, we also estimate the mean function implied by the restricted quartic model given by equation (18) in Murphy and Welch (1990):

$$\rho(\gamma, \theta_1, \theta_2, \theta_3, \alpha) = \theta_1 + \theta_2\gamma + (\theta_3 + \theta_2\alpha)\gamma^2 + 2\alpha\theta_3\gamma^3 + \alpha^2\theta_3\gamma^4 \quad (11)$$

that is commonly estimated in the empirical literature. Note that (11) is motivated as a more parsimonious model than the quartic model, and this formulation benefits from lessons in the empirical literature. If the model in (11) is misspecified, Theorem 1 indicates its estimated FMSE will be greater than that of a correctly specified model. Henceforth we denote this model as the quartic(r) model.

<<<<<< Insert Table 6.5. >>>>>>

The top panel of Table 6.5 reports the estimated FMSEs obtained by the quadratic, cubic, quartic, and restricted quartic models in parallel to Tables 2, 6, and 8 in Murphy and Welch (1990). As expected, the quartic specification provides the smallest FMSE and the quadratic specification yields the largest FMSE among the three specifications. We further observe that the FMSE substantially drops as the degree of the polynomial model increases, matching the improved fit of the MLIPs. The mean paths implied by the cubic and quartic functions are found to be statistically distinct. In particular, the FMSEs of both quartic and quartic(r) are almost identical for the workers without degree, implying that the quartic(r) is a very likely form for the workers' income path, although it is not effective for the other workers as the FMSEs from the quartic and quartic(r) models are distinctively different.

The lower panel of Table 6.5 reports the FMSEs using the data obtained by scaling the original income paths. That is, these functional observations are obtained by dividing each individual original LIP with by integral of the corresponding LIP over the entire working lifetime. The MLIPs are obtained using the same methodology as before for each group, and the corresponding FMSEs are shown in the lower panel of the Table. This normalization helps to exclude possible absolute effects of the level of lifetime income on inference concerning the MLIPs. Nonetheless, as the results in Table 6.5 show, the outcomes are qualitatively the same as for the top panel, so that the higher the degree of the polynomial model, the better the approximation. Scaling does not affect this outcome.

Although unreported here, we tested model adequacy by comparing the quadratic versus cubic models and the cubic versus quartic models. This procedure is motivated by the sequential testing procedure in Cho and Phillips (2018b) so that an adequate model can be found in a parsimonious manner. The QLR test statistic was employed in these comparisons using the FMSEs in Table 6.5 and the null limit distributions obtained in Section 4.3. The test trivially rejected the null hypothesis for every group classified by gender and education levels. Given the statistically substantial drops of the FMSEs as the degree of polynomial model increases, this rejection of the null hypothesis is expected and again reveals that higher-degree polynomial functions better approximate the MLIPs than lower-degree functions.

We next examine the gender and education effects on the MLIP. For this purpose, we extend the previous model structures using a dummy variable for the gender of each individual. Thus, for each education level, the model for the curves becomes

$$\rho_i(\gamma, \theta_1^F, \theta_2^F, \theta_3^F, \theta_1^M, \theta_2^M, \theta_3^M) = (\theta_1^F + \theta_2^F\gamma + \theta_3^F\gamma^2)d_i + (\theta_1^M + \theta_2^M\gamma + \theta_3^M\gamma^2)(1 - d_i) \quad (12)$$

extending the previous quadratic model using  $d_i = 1$  for female gender and  $d_i = 0$  otherwise. We estimate the unknown parameters by FLS and test whether the corresponding parameters are equal across genders using the Wald, LM and QLR test



statistics. So the hypotheses of interest are

$$\mathbb{H}_o : \theta_{1*}^F = \theta_{1*}^M, \theta_{2*}^F = \theta_{2*}^M, \text{ and } \theta_{3*}^F = \theta_{3*}^M, \text{ versus } \mathbb{H}_a : \theta_{1*}^F \neq \theta_{1*}^M, \theta_{2*}^F \neq \theta_{2*}^M, \text{ or } \theta_{3*}^F \neq \theta_{3*}^M$$

and failure to reject  $\mathbb{H}_o$  provides evidence that the MLIPs do not differ between genders. Similar extensions were made and tests conducted for the cubic, quartic, and quartic(r) models.

<<<<<< Insert Table 6.6. >>>>>>

The results are shown in Table 6.6, which provides empirical evidence for gender effects on the MLIPs. The top panel of Table 6.6 reports test results using the original LIP samples. For the groups with college education (Bachelor, Master, Ph.D), the null hypothesis of mean equivalence is rejected at both 1% and 5% levels, indicating that the average income paths are very different across different genders. This finding is consistent with empirical results in the literature. For example, Wiswall and Zafar (2017) find evidence of a gender effect on income paths and associate this difference with differing job demands between different genders. More specifically, they suggest that job flexibility and job stability may be more important factors for women in job choice, whereas men place a relatively higher preference on earnings, thereby producing a gender effect. These findings are nuanced for the group without college education. There the difference in the mean functions is significant for the LM test statistic, but the QLR test does not reject the null hypothesis of mean equality at the 1% and 5% levels. At the same levels, the Wald test rejects the null for the cubic and quartic specifications but not for the quadratic specification. Notwithstanding this outcome, overall it is evident that the gender effect is dominant for each education level and model specification.

In the lower panel of Table 6.6, we report results for the scaled LIP samples. The overall results from the top panel of Table 6.6 are evident also in the scaled samples. Thus, the gender effect is evident for the MLIPs, although it is less significant with the QLR statistic in the no-college education group and the differences in the MLIPs are less significant in Master and Ph.D groups, suggesting that the MLIPs become less differentiated across gender as educational qualifications rise to these levels.

Next, we examine the education effect on the MLIP across different education levels within the same gender. A pairwise comparison is made between the groups without college education and with a Bachelor degree and similar pairwise comparisons are made for those between the Bachelor and Masters level education, and Masters and Doctoral degrees.

<<<<<< Insert Table 6.7. >>>>>>

Table 6.7 reports the associated test outcomes. In the top panel of Table 6.7 we report all the results from the original LIP samples. At the 1% and 5% significance levels, the differences in the MLIPs turn out to be highly significant across different education levels for both male and female workers, implying that education levels have a big impact on the MLIPs, just as Mincer's (1958, 1974) wage equation implies. Since education levels are captured by dummy variables, we cannot be assured of linearity in the income path with respect to education years as Mincer's (1958, 1974) wage equation implies. Nevertheless, our finding is consistent with that wage equation.

The lower panel of Table 6.7 reports results using the scaled LIP samples. The overall education effect is evident for the scaled LIP samples just as before. But some qualitatively different results from the original samples are also apparent. In particular, in the male group, the mean difference is significant in the comparison between groups with no-college education and a Bachelor degree, but this difference diminishes for comparisons between groups with higher education at Masters and Doctoral levels. Thus, as the education level rises, the effects of education on the MLIP diminish for male workers. On the other hand, the female

group comparisons always yield  $p$ -values close to zero in all cases, implying that the female worker LIP is strictly affected by education at all levels and undiminished at the higher levels compared with male workers.

In sum, gender and education effects on the MLIP are evident in the LIP samples over the entire working lifetime. Different genders and education levels yield different MLIPs irrespective of whether each individual's LIP sample is scaled by aggregating over the full working cycle. Although there are some nuanced differences for different genders and education levels, the overall gender and education effects are evident in both cases.

## 6.2 Inference on the Mean of Log Income Path on the Post Work Experience Years

In this section we estimate the MLIP curves over the mature work experience years using the quadratic, cubic, quartic, and quartic(r) models as in Section 6.1. When the LIP samples are collected from income profiles with low job mobility, different gender and education effects are to be expected for the MLIP, as Mincer and Jovanovic (1981) point out.

<<<<<< Insert Figure 2. >>>>>>

Using the same samples and the same methodology as before, we estimate the shapes of the MLIPs on the work experience years from 10 to 40 years. For the same groups classified by gender and education levels, Figure 2 displays the estimated MLIPs implied by the quadratic, cubic, and quartic models along the pointwise estimated MLIP. As seen in Figure 2, the MLIPs implied by the cubic and quartic models are close to the overall pointwise MLIP. On the other hand, the quadratic model yields almost a linear MLIP over the mature work experience years and relatively large differences exist between the pointwise MLIP and the MLIP implied by the quadratic model.

<<<<<< Insert Table 6.8. >>>>>>

As before we proceed in parallel to Section 6.1 by first reporting in Table 6.8 the FMSEs obtained from the quadratic, cubic, quartic, and quartic(r) models. The format of Table 6.8 is the same as that of Table 6.5, but the estimated FMSE results exhibit different patterns from those in Table 6.5. For example, the top panel of Table 6.8 shows that the FMSE drops substantially as polynomial degree increases from quadratic to cubic for both genders and for all levels of education; but the FMSEs are more or less similar among the cubic, quartic, and quartic(r) models. So, extending model specification from cubic to quartic has little effect on model fit compared with increasing polynomial degree from quadratic to cubic. This outcome is concordant with Figure 2, in which the MLIP implied by the quadratic model differs from the pointwise MLIP and the MLIPs implied by the cubic and quartic models.

A similar pattern is observed in the bottom panel of Table 6.8, where FMSEs obtained from the scaled LIP samples are obtained. The scaled LIP paths are constructed differently from those in Section 6.1. Instead of dividing each individual's LIP with integrated LIP over lifetime work experience years, individual LIP is scaled by integrated LIP over mature work experience years. But irrespective of whether or not the LIPs are scaled, the same pattern is observed for the FMEs as in the top panel of Table 6.8.

Notwithstanding these results, it does not follow that higher than quadratic polynomial degrees are unhelpful in reducing the FMSEs. Although we do not report the results here, QLR tests of the quadratic specification versus the cubic specification and cubic versus quartic were conducted and these tests all rejected the null as before, implying that higher degree specifications deliver statistically more satisfactory approximations for the MLIPs.

<<<<<< Insert Table 6.9. >>>>>>

We next examine gender effects on the MLIPs over mature work experience years. As before, we apply Wald, LM, and QLR test statistics and report test outcomes in the top panel of Table 6.9 in parallel with Table 6.6. As seen in the top panel of Table 6.9, gender effects on the MLIP are evident in the groups with college educations. In particular, with each of the four model specifications, the hypothesis of constant MLIP between gender is rejected for groups with college education. In contrast, for the no-college education group it is difficult to reject constancy of the MLIP curves between gender, implying that female and male workers with low education levels face almost the same income profiles.

In the lower panel of Table 6.9, we report results using the scaled LIP samples. The figures in the table are obtained by conducting the same inferential procedures as those for the top panel. As is apparent in the table, no strong statistical evidence emerges indicating gender effects on the scaled MLIPs. Thus, if the LIP is associated with low job mobility, individual LIPs differ proportionately between gender so that gender influences are eliminated if proportionality is properly accounted for. This finding differs from that obtained from the LIPs over lifetime work experience years. If each individual's LIP is scaled by integrated LIP over the lifetime work years, the scaled LIP *is* affected by income profiles with high job mobility, leading to different MLIPs between different genders even when the LIPs are restricted to the mature work experience years.

<<<<<< Insert Table 6.10. >>>>>>

We close these empirics with an exploration of the effects of education within the same gender. Pairwise comparison results are reported in Table 6.10 in the same manner as before. The top panel gives test outcomes in parallel to those of Table 6.7. Using the original LIP samples, all test statistics produce qualitatively the same conclusions as those the top panel of Table 6.7, so that education effects on the MLIP are found to be present even after income profiles are constructed from workers with low job mobility.

However, when individual LIPs are scaled by integrated mature work experience levels, different results are obtained. The test outcomes are shown in the lower panel of Table 6.10. Just as in Table 6.9, no strong statistical evidence is found supporting different MLIPs among different education levels, giving the same conclusion as for the gender effect so that for the scaled LIPs education effects disappear within the same gender. Thus, if income profiles are generated from workers with low job mobility, individual LIPs differ proportionally among different education levels. So when these proportional differences are properly accounted for in the LIP samples, the MLIPs match across education levels.

In summary, if the LIP samples are constructed from workers with low job mobility, the empirical analysis produces different conclusions depending on whether the LIP is scaled or not. If the LIP samples are unscaled, results are overall similar to those obtained from the LIPs over the entire lifetime work experience years, so that the MLIPs obtained from these LIP samples do differ between gender and across education levels. But when the LIPs are scaled by respective integrated mature work experience years, the estimated MLIPs match, thereby implying that different gender and education levels affect the LIPs proportionally without changing the overall shape of the MLIP.

Before concluding, we provide some additional comments on the empirical analysis. First, all empirical models were also estimated using a beta distribution for the adjunct probability measure and estimating the shape parameters. This extension led to very similar estimation results and did not qualitatively change any of the findings. Use of the simpler uniform measure therefore seems adequate for this empirical application and, together with the simulations earlier discussed, implies that potential efficiency gains from this extension of FLS estimation may not be substantial. Second, the current empirical example assumes that  $x_i$  is discrete, but this is not needed for empirical analysis. Even when  $x_i$  is continuous, we can apply our methodology by forming  $\mathcal{M}$  appropriately and estimating the unknown parameter by FLS. For example, different education years can lead to different

LIPs, so that the coefficients in the quartic model can be modified to be a function affected by education years and relevant parameters to estimate. Finally, although the results are not reported here, all the empirical models were also estimated by first obtaining functional observations via the moving window method. This additional empirical analysis was conducted to remove any temporary shocks from the income profiles, but led to the same general conclusions reported above.

## 7 Conclusion

We develop a methodology of estimation and inference for parametric conditional mean functions that involve functional data. The approach is based on functional least squares (FLS) regression. Consistency and asymptotic normality are established under regularity conditions that allow for the possible presence of nuisance parameters and consequential effects on the limit theory. New Wald, Lagrange multiplier (LM), and quasi-likelihood ratio (QLR) tests are developed for this functional data context that enable inference about curve shapes in the observed data. Various examples where this methodology is useful include inference about the time forms in panel data observations, random coefficient model formalations, and distributional specifications in mixtures and copulas. The empirical functional form of income paths over work experience years is studied using these methods to examine the mean log income path within the framework of Mincer's (1958, 1974) wage equation. The findings show that gender and education levels produce differences in mean income paths but that these mean paths are proportional to each other. It follows that, upon rescaling the income paths by integrated work experience years associated with low job mobility, the mean income paths match over gender and across education levels.

The present study suggests potentially fruitful further developments in functional data inference. This paper assumes that explanatory variables are observed as vector valued whereas the dependent variable is observed in the form of random functions. Other cases of interest arise where both dependent and explanatory variables may take function space form or where the explanatory variable involves functional data and the dependent variable is scalar valued. See the model analyses in James, Wang, and Zhu (2009), Crambes, Gannoun, and Henchiri (2013), Petersen and Müller (2016), and Happ and Greven (2018) for some recent examples. Also, when different estimation methods than FLS are employed, new possibilities for parametric inference become possible. For instance, using penalized profile likelihood function estimation with functional data makes it possible to identify latent structures in functional curve data in a similar fashion to the classification Lasso method of Su, Shi, and Phillips (2016).

Further extensions of the present methods can be made to dependent data. Recently, functional time series analysis has attracted much attention in the literature. For instance, functional data such as the distribution of cross-sectional earnings or near-continuous recording of intraday stock returns potentially lead to autocorrelated dependence structures. Various methods of using such data can be found in Chang, Kim, and Park (2016), Kim and Park (2017), Beare, Seo, and Seo (2017), Hörmann, Kokoszka, and Nisol (2018), Seo and Beare (2019), Li, Robinson, and Shang (2020), Franchi and Paruolo (2020), and Chang, Hu, and Park (2019). From the perspective of parametric conditional mean function estimation, the presence of temporal dependence inevitably leads to limit theory for FLS estimation and statistical tests that differs from the present study and needs to be accounted for in time series inference. These are some of the many aspects of conditional mean function inference that can be addressed in future research.

# Appendices

These Appendices are organized as follows. Preliminary lemmas and proofs are given in Appendix A. Proofs of the results in the main text are given in Appendix B. The final Appendix C provides additional discussion and results on the estimation of the unconditional mean function.

## A Preliminary Lemmas

Before proving the claims in the text we provide some supplementary lemmas to be used later.

**Lemma 2.** *Given a measurable function  $h(\cdot, \theta) : \Gamma \mapsto \mathbb{R}$  on  $(\Gamma, \mathcal{G}, \mathbb{Q})$  for each  $\theta \in \Theta$ , if for each  $\gamma \in \Gamma$ ,  $h(\gamma, \cdot) \in \mathcal{C}^{(1)}(\Theta)$  and for each  $j \in \{1, 2, \dots, d\}$ ,  $\sup_{\theta \in \Theta} |(\partial/\partial\theta_j)h(\cdot, \theta)| \in L^1(\mathbb{Q})$ , then*

$$\nabla_{\theta} \int h(\gamma, \theta) d\mathbb{Q}(\gamma) = \int \nabla_{\theta} h(\gamma, \theta) d\mathbb{Q}(\gamma), \quad (13)$$

where  $\Theta$  is a compact and convex set in  $\mathbb{R}^d$  and  $d \in \mathbb{N}$  as in the text. □

**Lemma 3.** *Given Assumptions 1, 2, 3, and 4, for each  $\theta \in \Theta$ ,*

- (i)  $\nabla_{\theta} q(\theta) = -2 \int \int \{\mu(\gamma, x) - \rho(\gamma, \theta, x)\} \nabla_{\theta} \rho(\gamma, \theta, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma);$
- (ii)  $\nabla_{\theta}^2 q(\theta) = 2 \int \int \{\nabla_{\theta} \rho(\gamma, \theta, x) \nabla'_{\theta} \rho(\gamma, \theta, x) - \{\mu(\gamma, x) - \rho(\gamma, \theta, x)\} \nabla_{\theta}^2 \rho(\gamma, \theta, x)\} d\mathbb{P}(x) d\mathbb{Q}(\gamma);$
- (iii)  $\nabla_{\theta} q_n(\theta) = -2n^{-1} \sum_{i=1}^n \int \{g_i(\gamma) - \rho_i(\gamma, \theta)\} \nabla_{\theta} \rho_i(\gamma, \theta) d\mathbb{Q}(\gamma)$  a.s.  $-\mathbb{P}$ ; and
- (iv)  $\nabla_{\theta}^2 q_n(\theta) = 2n^{-1} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \theta) \nabla'_{\theta} \rho_i(\gamma, \theta) - \{g_i(\gamma) - \rho_i(\gamma, \theta)\} \nabla_{\theta}^2 \rho_i(\gamma, \theta) d\mathbb{Q}(\gamma)$  a.s.  $-\mathbb{P}$ . □

**Lemma 4.** *Given Assumptions 1, 2, 3, and 5, then  $B = \tilde{B}$ , where*

$$\tilde{B} := \int \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \varepsilon(\gamma, \theta_*) \varepsilon(\tilde{\gamma}, \theta_*) \nabla'_{\theta} \rho(\gamma, \theta_*, x) d\mathbb{P}(g(\gamma), g(\tilde{\gamma}), x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma})$$

and  $\varepsilon(\gamma, \theta) := g(\gamma) - \rho(\gamma, \theta, x)$ . □

**Lemma 5.** *Given Assumptions 2, 6, 7, and 8, for each  $\theta \in \Theta$ ,*

- (i)  $\nabla_{\theta} q(\theta) = -2 \int \int \{\mu(\gamma, x) - \rho(\gamma, \theta, x)\} \nabla_{\theta} \rho(\gamma, \theta, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma);$
- (ii)  $\nabla_{\theta}^2 q(\theta) = 2 \int \int \{\nabla_{\theta} \rho(\gamma, \theta, x) \nabla'_{\theta} \rho(\gamma, \theta, x) - \{\mu(\gamma, x) - \rho(\gamma, \theta, x)\} \nabla_{\theta}^2 \rho(\gamma, \theta, x)\} d\mathbb{P}(x) d\mathbb{Q}(\gamma);$
- (iii)  $\nabla_{\theta} \hat{q}_n(\theta) = -2n^{-1} \sum_{i=1}^n \int \{\hat{g}_i(\gamma) - \rho_i(\gamma, \theta)\} \nabla_{\theta} \rho_i(\gamma, \theta) d\mathbb{Q}(\gamma)$  a.s.  $-\mathbb{P}$ ; and
- (iv)  $\nabla_{\theta}^2 \hat{q}_n(\theta) = 2n^{-1} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \theta) \nabla'_{\theta} \rho_i(\gamma, \theta) - \{\hat{g}_i(\gamma) - \rho_i(\gamma, \theta)\} \nabla_{\theta}^2 \rho_i(\gamma, \theta) d\mathbb{Q}(\gamma)$  a.s.  $-\mathbb{P}$ . □

**Proof of Lemma 2:** From the differentiability condition, the given function is Lipschitz continuous, so that  $|h(\gamma, \theta) - h(\gamma, \theta')| \leq m(\gamma) \|\theta - \theta'\|$ , where for each  $\gamma \in \Gamma$ , we let  $m(\gamma) := \sup_{j \in \{1, 2, \dots, d\}} \sup_{\theta \in \Theta} |(\partial/\partial\theta_j)h(\gamma, \theta)|$ . Therefore, the following bound holds

$$\frac{1}{\|\theta - \theta'\|} \left| \int h(\gamma, \theta') d\mathbb{Q}(\gamma) - \int h(\gamma, \theta) d\mathbb{Q}(\gamma) \right| \leq \int m(\gamma) d\mathbb{Q}(\gamma) < \infty,$$

which further implies that

$$\lim_{\theta' \rightarrow \theta} \frac{1}{\|\theta - \theta'\|} \left[ \int h(\gamma, \theta') d\mathbb{Q}(\gamma) - \int h(\gamma, \theta) d\mathbb{Q}(\gamma) \right] = \int \lim_{\theta' \rightarrow \theta} \frac{1}{\|\theta - \theta'\|} [h(\gamma, \theta') - h(\gamma, \theta)] d\mathbb{Q}(\gamma),$$

by the dominated convergence theorem (DCT). The left and right sides of this equality are respectively identical to the left and right sides of (13). This completes the proof.  $\blacksquare$

**Proof of Lemma 3:** (i) The left side is expanded as

$$\nabla_{\theta} \int \int \{g(\gamma) - \rho(\gamma, \theta, x)\}^2 d\mathbb{P}(g(\gamma), x) d\mathbb{Q}(\gamma) = \nabla_{\theta} \int \int \{-2\mu(\gamma, x)\rho(\gamma, \theta, x) + \rho^2(\gamma, \theta, x)\} d\mathbb{P}(x) d\mathbb{Q}(\gamma),$$

using the fact that  $\mu(\gamma, x) = \int g(\gamma) d\mathbb{P}(g(\gamma)|x)$ . Given this,  $\mu_i(\cdot) \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$  and for each  $j \in \{1, 2, \dots, d\}$ ,  $\sup_{\theta} |(\partial/\partial\theta_j)\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$  because  $|g_i(\cdot)| \leq m_i$  and  $\sup_{(\gamma, \theta)} |(\partial/\partial\theta_j)\rho_i(\gamma, \theta)| \leq m_i$  by Assumptions 3, so that  $\sup_{\theta \in \Theta} |\mu_i(\cdot)(\partial/\partial\theta_j)\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  by the Cauchy-Schwarz's inequality. Therefore, applying Lemma 2 yields

$$\nabla_{\theta} \int \int \mu(\gamma, x)\rho(\gamma, \theta, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) = \int \int \mu(\gamma, x)\nabla_{\theta}\rho(\gamma, \theta, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma). \quad (14)$$

Furthermore,  $\sup_{(\gamma, \theta)} |\rho_i(\gamma, \theta)| \leq m_i$  by Assumption 3, so that  $\sup_{\theta} |\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$ , implying that  $\sup_{\theta \in \Theta} |\rho_i(\cdot, \theta)(\partial/\partial\theta_j)\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$ , again by Cauchy-Schwarz. Applying Lemma 2 entails that

$$\nabla_{\theta} \int \int \rho^2(\gamma, \theta, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) = 2 \int \int \rho(\gamma, \theta, x)\nabla_{\theta}\rho(\gamma, \theta, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma). \quad (15)$$

The desired result follows by combining (14) and (15).

(ii) Note that  $\mu_i(\cdot) \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$ ,  $\sup_{\theta \in \Theta} |\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$  as shown in (i). Furthermore, for each  $j, j' \in \{1, 2, \dots, d\}$ ,  $\sup_{\theta \in \Theta} |(\partial/\partial\theta_j)\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$  and  $\sup_{\theta \in \Theta} |(\partial^2/\partial\theta_j\partial\theta_{j'})\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$  by Assumption 3. This implies that

$$\sup_{\theta \in \Theta} |(\partial/\partial\theta_j)\rho_i(\cdot, \theta)(\partial/\partial\theta_{j'})\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q}) \text{ a.s. } -\mathbb{P},$$

$$\sup_{\theta \in \Theta} |\mu(\cdot)(\partial^2/\partial\theta_j\partial\theta_{j'})\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q}) \text{ a.s. } -\mathbb{P}, \text{ and } \sup_{\theta \in \Theta} |\rho_i(\cdot, \theta)(\partial^2/\partial\theta_j\partial\theta_{j'})\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q}) \text{ a.s. } -\mathbb{P}$$

by Cauchy-Schwarz. Applying Lemma 2 leads to the desired result as for (i).

(iii) The left side is expanded as follows:

$$\nabla_{\theta} \frac{1}{n} \sum_{i=1}^n \int \{g_i(\gamma) - \rho_i(\gamma, \theta)\}^2 d\mathbb{Q}(\gamma) = \nabla_{\theta} \frac{1}{n} \sum_{i=1}^n \int \{-2g_i(\gamma)\rho_i(\gamma, \theta) + \rho_i^2(\gamma, \theta)\} d\mathbb{Q}(\gamma).$$

Furthermore, for each  $j$ , the Cauchy-Schwarz inequality leads to  $\sup_{\theta \in \Theta} |g_i(\cdot)(\partial/\partial\theta_j)\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  because  $\sup_{\theta} |(\partial/\partial\theta_j)\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q})$  and  $\sup_{\theta} |g_i(\gamma)| \leq m_i \in L^2(\mathbb{P})$  as shown in (i and ii) using Assumption 3. Hence, applying Lemma 2 leads to

$$\nabla_{\theta} \frac{1}{n} \sum_{i=1}^n \int g_i(\gamma)\rho_i(\gamma, \theta) d\mathbb{Q}(\gamma) = \frac{1}{n} \sum_{i=1}^n \int g_i(\gamma)\nabla_{\theta}\rho_i(\gamma, \theta) d\mathbb{Q}(\gamma) \quad (16)$$

a.s.  $-\mathbb{P}$ . Finally, combining (15) and (16) yields the desired result.

(iv) Given the left side, for each  $j, j' \in \{1, 2, \dots, d\}$ ,  $\sup_{\theta \in \Theta} |(\partial/\partial\theta_j)\rho_i(\cdot, \theta)(\partial/\partial\theta_{j'})\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  by Cauchy-Schwarz and  $\sup_{\theta \in \Theta} |(\partial/\partial\theta_j)\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q})$  as shown in (ii). Further, for each  $j, j' \in \{1, 2, \dots, d\}$ ,  $\sup_{\theta \in \Theta} \{|g_i(\cdot) - \rho_i(\cdot, \theta)\}(\partial^2/\partial\theta_j\partial\theta_{j'})\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  again by applying Cauchy-Schwarz because  $g_i(\cdot) \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$ ,

$$\sup_{\theta} |\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q}) \text{ a.s. } -\mathbb{P}, \text{ and } \sup_{\theta \in \Theta} |(\partial^2/\partial\theta_j\partial\theta_{j'})\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q}) \text{ a.s. } -\mathbb{P}$$

as shown in (i, ii, and iii). Therefore,

$$\begin{aligned} & \nabla_{\theta} \frac{1}{n} \sum_{i=1}^n \int \{g_i(\gamma) - \rho_i(\gamma, \theta)\} \nabla_{\theta} \rho_i(\gamma, \theta) d\mathbb{Q}(\gamma) \\ &= \frac{1}{n} \sum_{i=1}^n \int \{g_i(\gamma) - \rho_i(\gamma, \theta)\} \nabla_{\theta}^2 \rho_i(\gamma, \theta) - \nabla_{\theta} \rho_i(\gamma, \theta) \nabla'_{\theta} \rho_i(\gamma, \theta) d\mathbb{Q}(\gamma) \quad \text{a.s.} - \mathbb{P} \end{aligned}$$

by applying Lemma 2. This completes the proof.  $\blacksquare$

**Proof of Lemma 4:** We first focus to the internal integral with respect to  $\mathbb{P}$ . Note that

$$\begin{aligned} & \int \varepsilon(\gamma, \theta_*) \varepsilon(\tilde{\gamma}, \theta_*) d\mathbb{P}(g(\gamma), g(\tilde{\gamma})|x) = \int \{g(\gamma) - \rho(\gamma, \theta_*, x)\} \{g(\tilde{\gamma}) - \rho(\tilde{\gamma}, \theta_*, x)\} d\mathbb{P}(g(\gamma), g(\tilde{\gamma})|x) \\ &= \int \{g(\gamma) - \mu(\gamma, x)\} \{g(\tilde{\gamma}) - \mu(\tilde{\gamma}, x)\} d\mathbb{P}(g(\gamma), g(\tilde{\gamma})|x) + \{\mu(\gamma, x) - \rho(\gamma, \theta_*, x)\} \{\mu(\tilde{\gamma}, x) - \rho(\tilde{\gamma}, \theta_*, x)\} \\ &= \kappa(\gamma, \tilde{\gamma}|x) + \{\mu(\gamma, x) - \rho(\gamma, \theta_*, x)\} \{\mu(\tilde{\gamma}, x) - \rho(\tilde{\gamma}, \theta_*, x)\} \end{aligned}$$

using the fact that  $\int \{g(\gamma) - \mu(\gamma, x)\} \{\mu(\tilde{\gamma}, x) - \rho(\tilde{\gamma}, \theta_*, x)\} d\mathbb{P}(g(\gamma), g(\tilde{\gamma})|x) = 0$ . Hence

$$\begin{aligned} \tilde{B} &= \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \kappa(\gamma, \tilde{\gamma}|x) \nabla'_{\theta} \rho(\tilde{\gamma}, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma}) \\ &\quad + \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \{\mu(\gamma, x) - \rho(\gamma, \theta_*, x)\} d\mathbb{Q}(\gamma) \int \{\mu(\tilde{\gamma}, x) - \rho(\tilde{\gamma}, \theta_*, x)\} \nabla'_{\theta} \rho(\tilde{\gamma}, \theta_*, x) d\mathbb{Q}(\tilde{\gamma}). \end{aligned}$$

Now note that by definition in Assumption 5(ii),  $B := \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \kappa(\gamma, \tilde{\gamma}|x) \nabla_{\theta} \rho(\tilde{\gamma}, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma})$  and  $\int \nabla_{\theta} \rho(\gamma, \theta_*, x) \{\mu(\gamma, x) - \rho(\gamma, \theta_*, x)\} d\mathbb{Q}(\gamma) = 0$  by the definition of  $\theta_*$  and Lemma 3(i). This completes the proof.  $\blacksquare$

**Proof of Lemma 5:** (i and ii) Assumptions 6 and 7 imply Assumptions 1 and 3. Therefore, the proofs of Lemma 3(i and ii) are sufficient for the proofs of Lemma 5(i and ii).

(iii) Note that  $\nabla_{\theta} \sum_{i=1}^n \int \{\hat{g}_i(\gamma) - \rho_i(\gamma, \theta)\}^2 d\mathbb{P} d\mathbb{Q}(\gamma) = \nabla_{\theta} \sum_{i=1}^n \int \{-2\hat{g}_i(\gamma) \rho_i(\gamma, \theta) + \rho_i^2(\gamma, \theta)\} d\mathbb{Q}(\gamma)$  a.s.  $-\mathbb{P}$ . Furthermore, for each  $j$ ,  $\sup_{\theta \in \Theta} |\hat{g}_i(\cdot) \cdot (\partial/\partial \theta_j) \rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  by Cauchy-Schwarz,

$$\sup_{\theta} |(\partial/\partial \theta_j) \rho_i(\cdot, \theta)| \in L^2(\mathbb{Q}),$$

as proved in the proof of Lemma 3, and  $\tilde{g}_i \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$  by Assumption 7(i). Hence, applying Lemma 2 leads to

$$\nabla_{\theta} \frac{1}{n} \sum_{i=1}^n \int \hat{g}_i(\gamma) \rho_i(\gamma, \theta) d\mathbb{Q}(\gamma) = \frac{1}{n} \sum_{i=1}^n \int \hat{g}_i(\gamma) \nabla_{\theta} \rho_i(\gamma, \theta) d\mathbb{Q}(\gamma) \quad (17)$$

a.s.  $-\mathbb{P}$ . Finally, combining (15) and (17) yields the desired result.

(iv) Given the left side, for each  $j$  and  $j'$ ,  $\sup_{\theta \in \Theta} |(\partial/\partial \theta_j) \rho_i(\cdot, \theta) \cdot (\partial/\partial \theta_{j'}) \rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  as proved in the proof of Lemma 3. Furthermore, for each  $j$  and  $j'$ ,  $\sup_{\theta \in \Theta} |\{\hat{g}_i(\cdot) - \rho_i(\cdot, \theta)\} (\partial^2/\partial \theta_j \partial \theta_{j'}) \rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$ , because  $\tilde{g}_i \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$  by Assumption 7(i), and  $\sup_{\theta \in \Theta} |\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q})$  a.s.  $-\mathbb{P}$  and  $\sup_{\theta \in \Theta} |(\partial^2/\partial \theta_j \partial \theta_{j'}) \rho_i(\cdot, \theta)| \in L^2(\mathbb{Q})$

a.s. –  $\mathbb{P}$  as shown in the proof of Lemma 3. Therefore,

$$\nabla_{\theta} \frac{1}{n} \sum_{i=1}^n \int \{\hat{g}_i(\gamma) - \rho_i(\gamma, \theta)\} \nabla_{\theta} \rho_i(\gamma, \theta) d\mathbb{Q}(\gamma) = \frac{1}{n} \sum_{i=1}^n \int \{\hat{g}_i(\gamma) - \rho_i(\gamma, \theta)\} \nabla_{\theta}^2 \rho_i(\gamma, \theta) - \nabla_{\theta} \rho_i(\gamma, \theta) \nabla'_{\theta} \rho_i(\gamma, \theta) d\mathbb{Q}(\gamma)$$

a.s. –  $\mathbb{P}$  by applying Lemma 2, thereby completing the proof.  $\blacksquare$

## B Proofs of the Main Results

**Proof of Theorem 1:** Note that

$$\begin{aligned} \int \int \{g(\gamma) - \rho(\gamma, \theta, x)\}^2 d\mathbb{P}(g(\gamma), x) d\mathbb{Q}(\gamma) &= \int \int \{g(\gamma) - \mu(\gamma, x)\}^2 d\mathbb{P}(g(\gamma), x) d\mathbb{Q}(\gamma) \\ &\quad - 2 \int \int \{g(\gamma) - \mu(\gamma, x)\} \{\rho(\gamma, \theta, x) - \mu(\gamma, x)\} d\mathbb{P}(g(\gamma), x) d\mathbb{Q}(\gamma) + \int \int \{\rho(\gamma, \theta, x) - \mu(\gamma, x)\}^2 d\mathbb{P}(x) d\mathbb{Q}(\gamma). \end{aligned}$$

Further,  $\int \int \{g(\gamma) - \mu(\gamma, x)\}^2 d\mathbb{P} d\mathbb{Q}(\gamma) = \int \text{var}_{\mathbb{P}}[g_i(\gamma)|x] d\mathbb{P}(x) d\mathbb{Q}$ , and  $\int \int \{g(\gamma) - \mu(\gamma, x)\} d\mathbb{P}(g(\gamma)|x) \{\mu(\gamma, x) - \rho(\gamma, \theta, x)\} d\mathbb{P}(x) d\mathbb{Q}(\gamma) = 0$  because  $\int \{g(\gamma) - \mu(\gamma, x)\} d\mathbb{P}(g(\gamma)|x) = 0$ . The desired result now follows.  $\blacksquare$

**Proof of Theorem 2:** (i) The result is obtained by applying the SULLN and DCT. Specifically, from the definitions of  $q_n(\cdot)$  and  $q(\cdot)$ , for each  $\theta$ ,

$$q_n(\theta) = \frac{1}{n} \sum_{i=1}^n \int g_i^2(\gamma) d\mathbb{Q}(\gamma) - \frac{2}{n} \sum_{i=1}^n \int (g_i(\gamma) \rho_i(\gamma, \theta)) d\mathbb{Q} + \frac{1}{n} \sum_{i=1}^n \int \rho_i^2(\gamma, \theta) d\mathbb{Q} \quad \text{and} \quad (18)$$

$$q(\theta) = \int E_{\mathbb{P}}[g_i^2(\gamma)] d\mathbb{Q}(\gamma) - 2 \int \mathbb{E}_{\mathbb{P}}[\mu_i(\gamma) \rho_i(\gamma, \theta)] d\mathbb{Q} + \int \mathbb{E}_{\mathbb{P}}[\rho_i^2(\gamma, \theta)] d\mathbb{Q}. \quad (19)$$

Here, we use the fact that  $\mathbb{E}_{\mathbb{P}}[g_i(\gamma) \rho_i(\gamma, \theta)] = \mathbb{E}_{\mathbb{P}}[\mathbb{E}_{\mathbb{P}}[g_i(\gamma)|x] \rho_i(\gamma, \theta)] = \mathbb{E}_{\mathbb{P}}[\mu_i(\gamma) \rho_i(\gamma, \theta)]$  in deriving (19), and we can interchange the integral and sample average operators in computing (18) by virtue of the DCT, as shown in the proof of Lemma 3.

We now examine the limit of each element in the right sides of (18) and (19). First, Assumption 3(i) implies that

$$\int \left| \frac{1}{n} \sum_{i=1}^n g_i^2(\gamma) - E_{\mathbb{P}}[g_i^2(\gamma)] \right| d\mathbb{Q}(\gamma) \leq \frac{1}{n} \sum_{i=1}^n m_i^2 + E_{\mathbb{P}}[m_i^2] < \infty \text{ a.s. – } \mathbb{P},$$

so that we can apply the DCT, giving

$$\int \left| \frac{1}{n} \sum_{i=1}^n g_i^2(\gamma) - E_{\mathbb{P}}[g_i^2(\gamma)] \right| d\mathbb{Q}(\gamma) \rightarrow 0 \text{ a.s. – } \mathbb{P}.$$

Next, Assumptions 3(i and ii) imply that  $\sup_{\theta \in \Theta} |\rho_i(\cdot, \theta) g_i(\cdot)| \leq m_i^2 \in L^1(\mathbb{P})$ , so that

$$\begin{aligned} \sup_{\theta \in \Theta} \int \left| \left( \frac{1}{n} \sum_{i=1}^n g_i(\gamma) \rho_i(\gamma, \theta) - \mathbb{E}[\mu_i(\gamma) \rho_i(\gamma, \theta)] \right) \right| d\mathbb{Q}(\gamma) \\ \leq \int \sup_{\theta \in \Theta} \left| \left( \frac{1}{n} \sum_{i=1}^n g_i(\gamma) \rho_i(\gamma, \theta) - \mathbb{E}[\mu_i(\gamma) \rho_i(\gamma, \theta)] \right) \right| d\mathbb{Q}(\gamma) \rightarrow 0 \end{aligned} \quad (20)$$



a.s. –  $\mathbb{P}$ , again by the DCT.

Third, from the fact that  $\sup_{\theta \in \Theta} |\rho_i(\cdot, \theta)| \in L^2(\mathbb{Q})$  a.s. –  $\mathbb{P}$ , as shown in the proof of Lemma 3,

$$\sup_{\theta \in \Theta} \int \left| \left( \frac{1}{n} \sum_{i=1}^n \rho_i^2(\gamma, \theta) - \mathbb{E}[\rho_i^2(\gamma, \theta)] \right) \right| d\mathbb{Q}(\gamma) \leq \int \sup_{\theta \in \Theta} \left| \left( \frac{1}{n} \sum_{i=1}^n \rho_i^2(\gamma, \theta) - \mathbb{E}[\rho_i^2(\gamma, \theta)] \right) \right| d\mathbb{Q}(\gamma) \rightarrow 0 \quad (21)$$

a.s. –  $\mathbb{P}$ .

Finally, combining the above three facts gives

$$\begin{aligned} \sup_{\theta \in \Theta} |q_n(\theta) - q(\theta)| &\leq \int \left| \frac{1}{n} \sum_{i=1}^n g_i^2(\gamma) - \mathbb{E}[g_i^2(\gamma)] \right| d\mathbb{Q}(\gamma) \\ &\quad + 2 \sup_{\theta \in \Theta} \left| \int \frac{1}{n} \sum_{i=1}^n g_i(\gamma) \rho_i(\gamma, \theta) - \mathbb{E}[\mu_i(\gamma) \rho_i(\gamma, \theta)] d\mathbb{Q}(\gamma) \right| \\ &\quad + \sup_{\theta \in \Theta} \int \left| \left( \frac{1}{n} \sum_{i=1}^n \rho_i^2(\gamma, \theta) - \mathbb{E}[\rho_i^2(\gamma, \theta)] \right) \right| d\mathbb{Q}(\gamma) \rightarrow 0 \text{ a.s. } - \mathbb{P} \end{aligned}$$

as desired.

(ii) This result follows from the definition of  $\hat{\theta}_n$  and Theorem 2(i), given the fact that for each  $\gamma \in \Gamma$ ,  $\rho_i(\gamma, \cdot)$  is in  $\mathcal{C}^{(2)}(\Theta)$  a.s. –  $\mathbb{P}$ . ■

**Proof of Theorem 3:** We first note that  $n^{-1} \sum_{i=1}^n \int \{g_i(\gamma) - \rho_i(\gamma, \hat{\theta}_n)\} \nabla_{\theta} \rho_i(\gamma, \hat{\theta}_n) d\mathbb{Q}(\gamma) = 0$  by Lemma 3(iii) and the definition of  $\hat{\theta}_n$ . We apply the mean-value theorem to the element in the integral by Lemma 3(iv), so that for some  $\bar{\theta}_n$  between  $\theta_*$  and  $\hat{\theta}_n$ , it follows that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \int \{g_i(\gamma) - \rho_i(\gamma, \hat{\theta}_n)\} \nabla_{\theta} \rho_i(\gamma, \hat{\theta}_n) d\mathbb{Q}(\gamma) &= \frac{1}{n} \sum_{i=1}^n \int \{g_i(\gamma) - \rho_i(\gamma, \theta_*)\} \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma) \\ &\quad + \frac{1}{n} \sum_{i=1}^n \int \{-\nabla_{\theta} \rho_i(\gamma, \bar{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \bar{\theta}_n) + [g_i(\gamma) - \rho_i(\gamma, \bar{\theta}_n)] \nabla_{\theta}^2 \rho_i(\gamma, \bar{\theta}_n)\} d\mathbb{Q}(\hat{\theta}_n - \theta_*), \end{aligned}$$

so that

$$A_n \sqrt{n} (\hat{\theta}_n - \theta_*) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \int [g_i(\gamma) - \rho_i(\gamma, \theta_*)] \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma), \quad (22)$$

where

$$A_n := \left\{ \frac{1}{n} \sum_{i=1}^n \int \left\{ \nabla_{\theta} \rho_i(\gamma, \bar{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \bar{\theta}_n) - [g_i(\gamma) - \rho_i(\gamma, \bar{\theta}_n)] \nabla_{\theta}^2 \rho_i(\gamma, \bar{\theta}_n) \right\} d\mathbb{Q}(\gamma) \right\}.$$

We now examine each element in (22), starting with the matrix  $A_n$  and writing

$$A_n = \frac{1}{n} \sum_{i=1}^n \left\{ \int \left\{ \nabla_{\theta} \rho_i(\gamma, \bar{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \bar{\theta}_n) \right\} d\mathbb{Q}(\gamma) - \int \left\{ [g_i(\gamma) - \rho_i(\gamma, \bar{\theta}_n)] \nabla_{\theta}^2 \rho_i(\gamma, \bar{\theta}_n) \right\} d\mathbb{Q}(\gamma) \right\}.$$

First, as already seen in the proof of Lemma 3(iv), for each  $j, j' \in \{1, 2, \dots, d\}$ ,  $|(\partial/\partial\theta_j) \rho_i(\cdot, \theta) \cdot (\partial/\partial\theta_{j'}) \rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s. –  $\mathbb{P}$ , so that Theorem 2 and interchanging the integral and sample average operators in the limit by the DCT yield

$$\int \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \rho_i(\gamma, \bar{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \bar{\theta}_n) d\mathbb{Q}(\gamma) \rightarrow \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \nabla'_{\theta} \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) \text{ a.s. } - \mathbb{P}.$$

Second, we note that for each  $i, j \in \{1, 2, \dots, d\}$ ,  $\sup_{\theta \in \Theta} |\rho_i(\cdot, \theta)(\partial^2/\partial\theta_j\partial\theta_{j'})\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  as shown in the proof of Lemma 3. Therefore, Theorem 2 and the DCT yield

$$\int \frac{1}{n} \sum_{i=1}^n \rho_i(\gamma, \bar{\theta}_n) \nabla_{\theta}^2 \rho_i(\gamma, \bar{\theta}_n) d\mathbb{Q}(\gamma) \rightarrow \int \int \rho(\gamma, \theta_*, x) \nabla_{\theta}^2 \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) \quad \text{a.s.} - \mathbb{P}.$$

Third, for each  $j, j' \in \{1, 2, \dots, d\}$ ,  $n^{-1} \sum_{i=1}^n \sup_{\theta \in \Theta} |g_i(\cdot) \cdot (\partial^2/\partial\theta_j\partial\theta_{j'})\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  as shown in the proof of Lemma 3. Therefore, in the same manner

$$\int \frac{1}{n} \sum_{i=1}^n g_i(\gamma) \nabla_{\theta}^2 \rho_i(\gamma, \bar{\theta}_n) d\mathbb{Q}(\gamma) \rightarrow \int \int \mu(\gamma, x) \nabla_{\theta}^2 \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) \quad \text{a.s.} - \mathbb{P}.$$

From these three facts, we deduce that  $A_n \rightarrow A$  a.s.  $-\mathbb{P}$ .

Next consider the right side of (22). By virtue of Assumption 1, we can apply the multivariate Lindeberg CLT giving

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \int [g_i(\gamma) - \rho_i(\gamma, \theta_*)] \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma) \stackrel{\Delta}{\sim} \mathcal{N}(0, B), \quad (23)$$

because the common mean of the components is  $\int \int \{g(\gamma) - \rho(\gamma, \theta_*, x)\} \nabla_{\theta} \rho(\gamma, \theta_*, x) d\mathbb{P}(g(\gamma), x) d\mathbb{Q}(\gamma) = \int \int \{\mu(\gamma, x) - \rho(\gamma, \theta_*, x)\} \nabla_{\theta} \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) = \nabla_{\theta} q(\theta_*) = 0$  by Lemma 3(i) and for each  $j$  and  $j'$ ,  $\int \int \int (\partial/\partial\theta_j)\rho(\gamma, \theta_*, x) \cdot \kappa(\gamma, \tilde{\gamma}|x) \cdot (\partial/\partial\theta_{j'})\rho(\tilde{\gamma}, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma}) < \infty$ , leading to a common positive definite variance matrix  $B$  by Assumption 5. From these properties and the IID condition of Assumption 1(ii), the right side of (22) is asymptotically distributed as  $\mathcal{N}(0, B)$ . Therefore,  $\sqrt{n}(\hat{\theta}_n - \theta_*) \stackrel{\Delta}{\sim} \mathcal{N}(0, A^{-1}BA^{-1})$ . This completes the proof.  $\blacksquare$

**Proof of Theorem 4:** (i) The desired result can be obtained by following the proof of Theorem 2. Specifically, we can apply the SULLN and DCT. From the definitions of  $\hat{q}_n(\cdot)$  and  $q(\cdot)$ , for each  $\theta$ ,

$$\hat{q}_n(\theta) = \int \frac{1}{n} \sum_{i=1}^n \hat{g}_i^2(\gamma) d\mathbb{Q}(\gamma) - 2 \int \frac{1}{n} \sum_{i=1}^n \{\hat{g}_i(\gamma) \rho_i(\gamma, \theta)\} d\mathbb{Q} + \int \frac{1}{n} \sum_{i=1}^n \rho_i^2(\gamma, \theta) d\mathbb{Q}, \quad (24)$$

interchanging the integral and finite sample averaging operators. We compare this expression with each element on the right side of (19). First, Assumption 7(i) implies that

$$\int \left| \frac{1}{n} \sum_{i=1}^n \hat{g}_i^2(\gamma) - E_{\mathbb{P}}[g_i^2(\gamma)] \right| d\mathbb{Q}(\gamma) \leq \frac{1}{n} \sum_{i=1}^n m_i^2 + E_{\mathbb{P}}[m_i^2] < \infty \quad \text{a.s.} - \mathbb{P},$$

so that application of the DCT gives

$$\int \left| \frac{1}{n} \sum_{i=1}^n \hat{g}_i^2(\gamma) - E_{\mathbb{P}}[g_i^2(\gamma)] \right| d\mathbb{Q}(\gamma) \rightarrow 0 \quad \text{a.s.} - \mathbb{P}$$

because: (a)

$$\int \left| \frac{1}{n} \sum_{i=1}^n \hat{g}_i^2(\gamma) - E_{\mathbb{P}}[g_i^2(\gamma)] \right| d\mathbb{Q}(\gamma) \leq \int \left| \frac{1}{n} \sum_{i=1}^n \hat{g}_i^2(\gamma) - \frac{1}{n} \sum_{i=1}^n g_i^2(\gamma) \right| d\mathbb{Q}(\gamma) + \int \left| \frac{1}{n} \sum_{i=1}^n g_i^2(\gamma) - E_{\mathbb{P}}[g_i^2(\gamma)] \right| d\mathbb{Q}(\gamma);$$

(b) the proof of Theorem 2(i) implies that  $\int \left| \frac{1}{n} \sum_{i=1}^n g_i^2(\gamma) - E_{\mathbb{P}}[g_i^2(\gamma)] \right| d\mathbb{Q}(\gamma) \rightarrow 0$  a.s.  $-\mathbb{P}$ ; and (c) by applying the mean-value

theorem for some  $\bar{\xi}_{n,\gamma}$  between  $\xi_*$  and  $\hat{\xi}_n$ ,

$$\int \left| \frac{1}{n} \sum_{i=1}^n \hat{g}_i^2(\gamma) - \frac{1}{n} \sum_{i=1}^n g_i^2(\gamma) \right| d\mathbb{Q}(\gamma) = 2 \int \left| \frac{1}{n} \sum_{i=1}^n \hat{g}_i(\gamma) \nabla_{\xi} \tilde{g}_i(\gamma, \bar{\xi}_{n,\gamma}) \right| d\mathbb{Q} \cdot |\hat{\xi}_n - \xi_*|.$$

Thus, for each  $j = 1, 2, \dots, s$ , Assumption 7(iii) implies that  $\int \left| \frac{1}{n} \sum_{i=1}^n \hat{g}_i(\gamma) \cdot (\partial/\partial \xi_j) g_i(\gamma, \bar{\xi}_{n,\gamma}) \right| d\mathbb{Q} \leq n^{-1} \sum_{i=1}^n m_i^2$ , and  $\hat{\xi}_n \rightarrow \xi_*$  a.s.  $-\mathbb{P}$  by Assumption 8(i). The desired result follows from properties (a), (b), and (c).

Next, compare the second elements in (24) and (19). First,

$$\sup_{\theta \in \Theta} \left| \int \frac{1}{n} \sum_{i=1}^n (\hat{g}_i(\gamma) \rho_i(\gamma, \theta) - \mathbb{E}_{\mathbb{P}}[\mu_i(\gamma) \rho_i(\gamma, \theta)]) d\mathbb{Q} \right| = \int \sup_{\theta \in \Theta} \left| \left( \frac{1}{n} \sum_{i=1}^n g_i(\gamma) \rho_i(\gamma, \theta) - \mathbb{E}_{\mathbb{P}}[g_i(\gamma) \rho_i(\gamma, \theta)] \right) \right| d\mathbb{Q} + o_{\mathbb{P}}(1),$$

because applying the mean-value theorem implies that for some  $\bar{\xi}_{\gamma,n}$  between  $\xi_*$  and  $\hat{\xi}_n$ ,

$$\sup_{\theta \in \Theta} \left| \int \left( \frac{1}{n} \sum_{i=1}^n \{\hat{g}_i(\gamma) - g_i(\gamma)\} \right) \rho_i(\gamma, \theta) d\mathbb{Q} \right| = \sup_{\theta \in \Theta} \left| \int \frac{1}{n} \sum_{i=1}^n \nabla_{\xi} \tilde{g}_i(\gamma, \bar{\xi}_{\gamma,n}) \cdot \rho_i(\gamma, \theta) d\mathbb{Q} \cdot (\hat{\xi}_n - \xi_*) \right|,$$

so that for each  $j = 1, 2, \dots, s$ ,

$$\sup_{\theta \in \Theta} \left| \int \frac{1}{n} \sum_{i=1}^n (\partial/\partial \xi_j) \tilde{g}_i(\gamma, \xi_{\gamma,n}) \cdot \rho_i(\gamma, \theta) d\mathbb{Q} \cdot (\hat{\xi}_{j,n} - \xi_{j*}) \right| \leq \left( \frac{1}{n} \sum_{i=1}^n m_i^2 \right)^{1/2} \sup_{\theta \in \Theta} \left| \int \frac{1}{n} \sum_{i=1}^n \rho_i^2(\gamma, \theta) d\mathbb{Q} \right| \cdot |\hat{\xi}_{j,n} - \xi_{j*}| \rightarrow 0$$

a.s.  $-\mathbb{P}$  by Assumptions 7(ii and iii) and 8(i). Now, (21) implies that

$$\sup_{\theta \in \Theta} \left| \int \frac{1}{n} \sum_{i=1}^n (\hat{g}_i(\gamma) \rho_i(\gamma, \theta) - \mathbb{E}_{\mathbb{P}}[\mu_i(\gamma) \rho_i(\gamma, \theta)]) d\mathbb{Q} \right| \rightarrow 0 \text{ a.s. } -\mathbb{P}.$$

Finally, the third component in the right side of (24) is identical to the third element in the right side of (19). The desired result follows from these three properties.

(ii) This result follows from the definition of  $\tilde{\theta}_n$  and Theorem 4(i), given the fact that for each  $\gamma \in \Gamma$ ,  $\rho_i(\gamma, \cdot)$  is in  $\mathcal{C}^{(2)}(\Theta)$  a.s.  $-\mathbb{P}$ . ■

**Proof of Theorem 5:** We first note that  $n^{-1} \sum_{i=1}^n \int \{\hat{g}_i(\gamma) - \rho_i(\gamma, \tilde{\theta}_n)\} \nabla_{\theta} \rho_i(\gamma, \tilde{\theta}_n) d\mathbb{Q}(\gamma) \equiv 0$  by Lemma 5(iii) and the definition of  $\tilde{\theta}_n$ . We apply the mean-value theorem to the element in the integral by Lemma 3(iv), so that for some  $\bar{\theta}_n$  between  $\theta_*$  and  $\tilde{\theta}_n$ , it follows that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \int \{\hat{g}_i(\gamma) - \rho_i(\gamma, \tilde{\theta}_n)\} \nabla_{\theta} \rho_i(\gamma, \tilde{\theta}_n) d\mathbb{Q}(\gamma) &= \frac{1}{n} \sum_{i=1}^n \int \{\hat{g}_i(\gamma) - \rho_i(\gamma, \theta_*)\} \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma) \\ &+ \frac{1}{n} \sum_{i=1}^n \int \{-\nabla_{\theta} \rho_i(\gamma, \bar{\theta}_n) \nabla_{\theta} \rho_i(\gamma, \bar{\theta}_n) + [g_i(\gamma) - \rho_i(\gamma, \bar{\theta}_n)] \nabla_{\theta}^2 \rho_i(\gamma, \bar{\theta}_n)\} d\mathbb{Q}(\bar{\theta}_n - \theta_*), \end{aligned}$$

and then

$$\hat{A}_n \sqrt{n} (\tilde{\theta}_n - \theta_*) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ \int [\hat{g}_i(\gamma) - \rho_i(\gamma, \theta_*)] \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma) \right\}, \quad (25)$$

where

$$\widehat{A}_n := \left\{ \int \{ \nabla_{\theta} \rho_i(\gamma, \bar{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \bar{\theta}_n) - \frac{1}{n} \sum_{i=1}^n [\widehat{g}_i(\gamma) - \rho_i(\gamma, \bar{\theta}_n)] \nabla_{\theta}^2 \rho_i(\gamma, \bar{\theta}_n) \} d\mathbb{Q}(\gamma) \right\}.$$

We examine each element in (25). First, for each  $j, j' \in \{1, 2, \dots, d\}$ ,  $|(\partial/\partial\theta_j)\rho_i(\cdot, \theta) \cdot (\partial/\partial\theta_{j'})\rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  as shown in the proof of Lemma 3, implying that

$$\frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \bar{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \bar{\theta}_n) d\mathbb{Q}(\gamma) \rightarrow \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \nabla'_{\theta} \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) \text{ a.s. } - \mathbb{P}.$$

by Theorem 4(ii) and the DCT.

Second, for each  $i, j \in \{1, 2, \dots, d\}$ ,  $\sup_{\theta \in \Theta} |\rho_i(\cdot, \theta) (\partial^2/\partial\theta_j \partial\theta_{j'}) \rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  by Cauchy-Schwarz, Assumption 7(i) and 4(i). Therefore,

$$\int \frac{1}{n} \sum_{i=1}^n \rho_i(\gamma, \bar{\theta}_n) \nabla_{\theta}^2 \rho_i(\gamma, \bar{\theta}_n) d\mathbb{Q}(\gamma) \rightarrow \int \int \rho(\gamma, \theta_*, x) \nabla_{\theta}^2 \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) \text{ a.s. } - \mathbb{P}$$

by the DCT and Theorem 4(ii).

Third, for each  $j, j' \in \{1, 2, \dots, d\}$ ,  $n^{-1} \sum_{i=1}^n \sup_{\xi \in \Xi} \sup_{\theta \in \Theta} |\widetilde{g}_i(\cdot, \cdot, \xi) \cdot (\partial^2/\partial\theta_j \partial\theta_{j'}) \rho_i(\cdot, \theta)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$  as shown in the proof of Lemma 3. Thus, Assumption 3 and the DCT imply that

$$\frac{1}{n} \int \sum_{i=1}^n G_{iy}(\gamma) \nabla_{\theta}^2 \rho_i(\gamma, \bar{\theta}_n) d\mathbb{Q}(\gamma) \rightarrow \int \int \mu(\gamma, x) \nabla_{\theta}^2 \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) \text{ a.s. } - \mathbb{P}.$$

These results give  $\widehat{A}_n \rightarrow A$  a.s.  $-\mathbb{P}$ .

Next, we examine the right side of (25). Applying the mean-value theorem, we obtain the representation given in (5), viz., that for some  $\bar{\xi}_{n,\gamma}$  between  $\widehat{\xi}_n$  and  $\xi_*$ ,

$$\begin{aligned} & \frac{1}{\sqrt{n}} \int \sum_{i=1}^n \{ \widehat{g}_i(\gamma) - \rho_i(\gamma, \theta_*) \} \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma) \\ &= \frac{1}{\sqrt{n}} \int \sum_{i=1}^n \{ g_i(\gamma, \xi_*) - \rho_i(\gamma, \theta_*) \} \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma) + \frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \theta_*) \cdot \nabla'_{\xi} g_i(\gamma, \bar{\xi}_{n,\gamma}) d\mathbb{Q} \cdot \sqrt{n}(\widehat{\xi}_n - \xi_*). \end{aligned} \quad (26)$$

From the proof of Theorem 3 we have

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \int [g_i(\gamma) - \rho_i(\gamma, \theta_*)] \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma) \stackrel{\Delta}{\sim} \mathcal{N}(0, B).$$

We also note that for  $j = 1, \dots, d$  and  $j' = 1, \dots, s$ , Assumptions 3(iii) and 7(iii) imply that  $\sup_{(\theta, \xi)} |(\partial/\partial\theta_j)\rho_i(\cdot, \theta) \cdot (\partial/\partial\xi_{j'})\widetilde{g}_i(\cdot, \xi)| \in L^1(\mathbb{Q})$  a.s.  $-\mathbb{P}$ , so that applying the DCT shows that

$$\frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \theta_*) \nabla'_{\xi} g_i(\gamma, \bar{\xi}_{n,\gamma}) d\mathbb{Q}(\gamma) \rightarrow M := \int \mathbb{E}_{\mathbb{P}} [\nabla_{\theta} \rho(\gamma, \theta_*, x_i) \nabla'_{\xi} g_i(\gamma, \xi_*)] d\mathbb{Q}(\gamma) \text{ a.s. } - \mathbb{P}.$$

In addition, if we combine Assumptions 8(ii and iii) and 9,

$$\sqrt{n}(\widehat{\xi}_n - \xi_*) = -H^{-1} \sqrt{n} s_{*n} + o_{\mathbb{P}}(1) \stackrel{\Delta}{\sim} \mathcal{N}(0, H^{-1} J H^{-1'}).$$

Therefore,

$$\frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \theta_*) \nabla'_{\xi} g_i(\gamma, \bar{\xi}_{n,\gamma}) d\mathbb{Q} \sqrt{n}(\hat{\xi}_n - \xi_*) \stackrel{\Delta}{\approx} \mathcal{N}(0, MH^{-1}JH^{-1}'M^{-1}).$$

Furthermore, the asymptotic covariance between  $n^{-1/2} \sum_{i=1}^n \int [g_i(\gamma) - \rho_i(\gamma, \theta_*)] \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma)$  and  $n^{-1/2} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \theta_*) \cdot \nabla'_{\xi} g_i(\gamma, \bar{\xi}_{n,\gamma}) d\mathbb{Q}(\hat{\xi}_n - \xi_*)$  is obtained as  $-MH^{-1}K$  by Assumption 9, implying that the asymptotic covariance matrix of  $n^{-1/2} \int \sum_{i=1}^n \{\hat{g}_i(\gamma) - \mu_i(\gamma)\} \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma)$  is given by  $B_* := B - MH^{-1}K - K'H^{-1}M' + MH^{-1}JH^{-1}'M'$ , which is positive definite by Assumption 9. We therefore obtain

$$\frac{1}{\sqrt{n}} \int \sum_{i=1}^n \{\hat{g}_i(\gamma) - \mu_i(\gamma)\} \nabla_{\theta} \rho_i(\gamma, \theta_*) d\mathbb{Q}(\gamma) \stackrel{\Delta}{\approx} \mathcal{N}(0, B_*). \quad (27)$$

Finally,  $A^{-1}$  exists by Assumption 4(ii), so that  $\sqrt{n}(\tilde{\theta}_n - \theta_*) \stackrel{\Delta}{\approx} \mathcal{N}(0, A^{-1}B_*A^{-1})$  by (25) and (27). This completes the proof.  $\blacksquare$

**Proof of Theorem 6:** We first consider the consistency of  $\hat{A}_n$ . Note that for each  $j$  and  $j' = 1, 2, \dots, d$ ,  $\sup_{(\gamma, \theta) \in \Gamma \times \Theta} |(\partial/\partial\theta_j)\rho_i(\gamma, \theta)(\partial/\partial\theta_{j'})\rho_i(\gamma, \theta)| \leq m_i \in L^2(\mathbb{P})$  by Assumption 3(iii), so that

$$\sup_{(\gamma, \theta) \in \Gamma \times \Theta} \left| \frac{1}{n} \sum_{i=1}^n (\nabla_{\theta} \rho_i(\gamma, \theta) \nabla'_{\theta} \rho_i(\gamma, \theta) - \mathbb{E}_{\mathbb{P}}[\nabla_{\theta} \rho_i(\gamma, \theta) \nabla'_{\theta} \rho_i(\gamma, \theta)]) \right| \rightarrow 0 \text{ a.s.} - \mathbb{P}$$

by the SULLN. Therefore, applying the DCT, it follows that

$$\frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \hat{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \hat{\theta}_n) d\mathbb{Q}(\gamma) \rightarrow \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \nabla'_{\theta} \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) \text{ a.s.} - \mathbb{P} \quad (28)$$

using the fact that  $\hat{\theta}_n \rightarrow \theta_*$  a.s.  $-\mathbb{P}$ . We also note that  $\sup_{(\gamma, \theta) \in \Gamma \times \Theta} |\varepsilon_i(\gamma, \theta)(\partial^2/\partial\theta_j\partial\theta_{j'})\rho_i(\gamma, \theta)| \leq m_i^2 \in L^1(\mathbb{P})$  by Assumptions 3(i, ii, and iv). Therefore,

$$\sup_{(\gamma, \theta) \in \Gamma \times \Theta} \left| \frac{1}{n} \sum_{i=1}^n (\varepsilon_i(\gamma, \theta) \cdot \nabla_{\theta}^2 \rho_i(\gamma, \theta) - \mathbb{E}_{\mathbb{P}}[\varepsilon_i(\gamma, \theta) \cdot \nabla_{\theta}^2 \rho_i(\gamma, \theta)]) \right| \rightarrow 0 \text{ a.s.} - \mathbb{P}$$

by the SULLN. Therefore, the DCT and the fact that  $\hat{\theta}_n \rightarrow \theta_*$  a.s.  $-\mathbb{P}$  imply that

$$\frac{1}{n} \sum_{i=1}^n \int \varepsilon_i(\gamma, \hat{\theta}_n) \cdot \nabla_{\theta}^2 \rho_i(\gamma, \hat{\theta}_n) d\mathbb{Q}(\gamma) \rightarrow \int \int \{\mu(\gamma, x) - \rho(\gamma, \theta_*, x)\} \nabla_{\theta}^2 \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) \text{ a.s.} - \mathbb{P}. \quad (29)$$

Here, we used the fact that  $\int g(\gamma) d\mathbb{P}(g(\gamma)|x) = \mu(\gamma, x)$ . Now, (28) and (29) imply that  $\hat{A}_n \rightarrow A$  a.s.  $-\mathbb{P}$ .

We next examine the consistency of  $\hat{B}_n$ . Note that for each  $j$  and  $j' = 1, 2, \dots, d$ ,  $\sup_{(\gamma, \tilde{\gamma}, \theta) \in \Gamma \times \Gamma \times \Theta} |(\partial/\partial\theta_j)\rho_i(\gamma, \theta)\varepsilon_i(\gamma, \theta)\varepsilon_i(\tilde{\gamma}, \theta)(\partial/\partial\theta_{j'})\rho_i(\tilde{\gamma}, \theta)| \leq m_i^2 \in L^1(\mathbb{P})$  by Assumptions 3(i, ii, and iii), so that

$$\sup_{(\gamma, \tilde{\gamma}, \theta) \in \Gamma \times \Gamma \times \Theta} \left| \frac{1}{n} \sum_{i=1}^n (\nabla_{\theta} \rho_i(\gamma, \theta) \varepsilon_i(\gamma, \theta) \varepsilon_i(\tilde{\gamma}, \theta) \nabla'_{\theta} \rho_i(\tilde{\gamma}, \theta) - \mathbb{E}_{\mathbb{P}}[\nabla_{\theta} \rho_i(\gamma, \theta) \varepsilon_i(\gamma, \theta) \varepsilon_i(\tilde{\gamma}, \theta) \nabla'_{\theta} \rho_i(\tilde{\gamma}, \theta)]) \right| \rightarrow 0$$

a.s. -  $\mathbb{P}$  by the SULLN. Therefore, applying the DCT, it follows that

$$\begin{aligned}\widehat{B}_n &:= \frac{1}{n} \sum_{i=1}^n \int \int \nabla_{\theta} \rho_i(\gamma, \widehat{\theta}_n) \varepsilon_i(\gamma, \widehat{\theta}_n) \varepsilon_i(\widetilde{\gamma}, \widehat{\theta}_n) \nabla'_{\theta} \rho_i(\widetilde{\gamma}, \widehat{\theta}_n) d\mathbb{Q}(\gamma) d\mathbb{Q}(\widetilde{\gamma}) \\ &\rightarrow \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \int \varepsilon(\gamma, \theta_*) \varepsilon(\widetilde{\gamma}, \theta_*) d\mathbb{P}(g(\gamma), g(\widetilde{\gamma})|x) \nabla'_{\theta} \rho(\widetilde{\gamma}, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\widetilde{\gamma}) \quad a.s. - \mathbb{P}\end{aligned}$$

using the fact that  $\widehat{\theta}_n \rightarrow \theta_*$  a.s. -  $\mathbb{P}$ . Here,

$$\int \varepsilon(\gamma, \theta_*) \varepsilon(\widetilde{\gamma}, \theta_*) d\mathbb{P}(g(\gamma), g(\widetilde{\gamma})|x) = \int \{g(\gamma) - \rho(\gamma, \theta_*, x)\} \{g(\widetilde{\gamma}) - \rho(\widetilde{\gamma}, \theta_*, x)\} d\mathbb{P}(|x) =: \kappa(\gamma, \widetilde{\gamma}|x).$$

Therefore,  $\widehat{B}_n \rightarrow \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \kappa(\gamma, \widetilde{\gamma}|x) \nabla'_{\theta} \rho(\widetilde{\gamma}, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\widetilde{\gamma})$  a.s. -  $\mathbb{P}$ , corresponding to the definition of  $B$ .

This completes the proof.  $\blacksquare$

**Proof of Theorem 7:** We start by showing consistency of  $\widetilde{A}_n$ . First, (28) implies that

$$\frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \widetilde{\theta}_n) \nabla'_{\theta} \rho_i(\gamma, \widetilde{\theta}_n) d\mathbb{Q}(\gamma) \rightarrow \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \nabla'_{\theta} \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma). \quad (30)$$

Next, Assumptions 11(i and iii) imply that  $\sup_{(\gamma, \theta, \xi) \in \Gamma \times \Theta \times \Xi} |\varepsilon_i(\gamma, \theta, \xi) \cdot (\partial^2 / \partial \theta_j \partial \theta_{j'}) \rho_i(\gamma, \theta)| \leq m_i^2 \in L^1(\mathbb{P})$  for each  $j$  and  $j' = 1, 2, \dots, d$ . Therefore, the SULLN implies that

$$\sup_{(\gamma, \theta, \xi) \in \Gamma \times \Theta \times \Xi} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i(\gamma, \theta, \xi) \cdot \nabla_{\theta}^2 \rho_i(\gamma, \theta) - \mathbb{E}_{\mathbb{P}}[\varepsilon_i(\gamma, \theta, \xi) \cdot \nabla_{\theta}^2 \rho_i(\gamma, \theta)] \right| \rightarrow 0 \quad a.s. - \mathbb{P}$$

by the DCT. Therefore,

$$\frac{1}{n} \sum_{i=1}^n \int \widetilde{\varepsilon}_i(\gamma, \widetilde{\theta}_n) \cdot \nabla_{\theta}^2 \rho_i(\gamma, \widetilde{\theta}_n) \rightarrow \int \int \{\mu(\gamma, x) - \rho(\gamma, \theta_*, x)\} \nabla_{\theta}^2 \rho(\gamma, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) \quad a.s. - \mathbb{P} \quad (31)$$

noting that  $\int g(\gamma) d\mathbb{P}(g(\gamma)|x) = \mu(\gamma, x)$ . Now, (30) and (31) imply that  $\widetilde{A}_n \rightarrow A$  a.s. -  $\mathbb{P}$ .

We next show the consistency of  $\widetilde{B}_n$ . From the definition of  $\widetilde{B}_n$ , if we show that (i)  $\widetilde{B}_n \rightarrow B$  a.s. -  $\mathbb{P}$ , (ii)  $\widehat{M}_n \rightarrow M$  a.s. -  $\mathbb{P}$ , and (iii)  $\widehat{K}_n \rightarrow K$  a.s. -  $\mathbb{P}$ , then the consistency of  $\widetilde{B}_n$  follows from Assumption 10.

(i) Proving that  $\widetilde{B}_n \rightarrow B$  a.s. -  $\mathbb{P}$  is almost identical to proving that of  $\widehat{B}_n \rightarrow B$ . Note that for each  $j$  and  $j' = 1, 2, \dots, d$ ,  $\sup_{(\gamma, \widetilde{\gamma}, \theta, \xi) \in \Gamma \times \Gamma \times \Theta} |(\partial / \partial \theta_j) \rho_i(\gamma, \theta) \varepsilon_i(\gamma, \theta, \xi) \varepsilon_i(\widetilde{\gamma}, \theta, \xi) (\partial / \partial \theta_{j'}) \rho_i(\widetilde{\gamma}, \theta)| \leq m_i^2 \in L^1(\mathbb{P})$  by Assumptions 11(i and ii), so that

$$\sup_{(\gamma, \widetilde{\gamma}, \theta, \xi)} \left| \frac{1}{n} \sum_{i=1}^n (\nabla_{\theta} \rho_i(\gamma, \theta) \varepsilon_i(\gamma, \theta, \xi) \varepsilon_i(\widetilde{\gamma}, \theta, \xi) \nabla'_{\theta} \rho_i(\widetilde{\gamma}, \theta) - \mathbb{E}_{\mathbb{P}}[\nabla_{\theta} \rho_i(\gamma, \theta) \varepsilon_i(\gamma, \theta, \xi) \varepsilon_i(\widetilde{\gamma}, \theta, \xi) \nabla'_{\theta} \rho_i(\widetilde{\gamma}, \theta)]) \right| \rightarrow 0$$

a.s. -  $\mathbb{P}$  by the SULLN. Therefore, applying the DCT, it follows that

$$\begin{aligned}\widetilde{B}_n &:= \frac{1}{n} \sum_{i=1}^n \int \int \nabla_{\theta} \rho_i(\gamma, \widetilde{\theta}_n) \widetilde{\varepsilon}_{in}(\gamma, \widetilde{\theta}_n) \widetilde{\varepsilon}_{in}(\widetilde{\gamma}, \widetilde{\theta}_n) \nabla'_{\theta} \rho_i(\widetilde{\gamma}, \widetilde{\theta}_n) d\mathbb{Q}(\gamma) d\mathbb{Q}(\widetilde{\gamma}) \\ &\rightarrow \int \int \nabla_{\theta} \rho(\gamma, \theta_*, x) \int \varepsilon(\gamma, \theta_*, \xi_*) \varepsilon(\widetilde{\gamma}, \theta_*, \xi_*) d\mathbb{P}(g(\gamma), g(\widetilde{\gamma})|x) \nabla'_{\theta} \rho(\widetilde{\gamma}, \theta_*, x) d\mathbb{P}(x) d\mathbb{Q}(\gamma) d\mathbb{Q}(\widetilde{\gamma}) \quad a.s. - \mathbb{P}\end{aligned}$$

using the fact that  $(\widehat{\xi}_n, \widetilde{\theta}_n) \rightarrow (\xi_*, \theta_*)$  a.s. -  $\mathbb{P}$ . In the proof of Theorem 6, we have already seen that the right side is identical to

B. Therefore,  $\bar{B}_n \rightarrow B$  a.s.  $-\mathbb{P}$ .

(ii) Now we show  $\widehat{M}_n \rightarrow M$  a.s.  $-\mathbb{P}$ . Note that Assumptions 11(ii and iv) imply that for each  $j = 1, 2, \dots, d$  and  $j' = 1, 2, \dots, s$ ,  $\sup_{(\gamma, \theta, \xi) \in \Gamma \times \Theta \times \Xi} |(\partial/\partial\theta_j)\rho_i(\gamma, \theta) \cdot (\partial/\partial\xi_{j'})\tilde{G}y_i(\gamma, \xi)| \leq m_i^2 \in L^1(\mathbb{P})$ . Therefore,

$$\sup_{(\gamma, \theta, \xi) \in \Gamma \times \Theta \times \Xi} \left| \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \rho_i(\gamma, \theta) \nabla'_{\xi} \tilde{g}_i(\gamma, \xi) - \mathbb{E}_{\mathbb{P}}[\nabla_{\theta} \rho_i(\gamma, \theta) \nabla'_{\xi} \tilde{g}_i(\gamma, \xi)] \right| \rightarrow 0 \text{ a.s. } -\mathbb{P} \quad (32)$$

by the SULLN. Therefore, applying the DCT implies that

$$\widehat{M}_n := \frac{1}{n} \sum_{i=1}^n \int \nabla_{\theta} \rho_i(\gamma, \tilde{\theta}_n) \nabla'_{\xi} \tilde{g}_i(\gamma, \hat{\xi}_n) d\mathbb{Q}(\gamma) \rightarrow \int E_{\mathbb{P}}[\nabla_{\theta} \rho_i(\gamma, \theta_*) \nabla'_{\xi} \tilde{g}_i(\gamma, \xi_*)] d\mathbb{Q}(\gamma) \text{ a.s. } -\mathbb{P}$$

using the fact that  $(\hat{\xi}_n, \tilde{\theta}_n) \rightarrow (\xi_*, \theta_*)$  a.s.  $-\mathbb{P}$ . Note that the right side is  $M$ , implying that  $\widehat{M}_n \rightarrow M$  a.s.  $-\mathbb{P}$ .

(iii) Next we show  $\widehat{K}_n \rightarrow K$  a.s.  $-\mathbb{P}$ . Note that  $\sup_{(\gamma, \theta, \xi) \in \Gamma \times \Theta \times \Xi} |s_i(\xi)\{\tilde{g}_i(\gamma, \xi) - \rho_i(\gamma, \theta)\} \nabla'_{\theta} \rho_i(\gamma, \theta)| \leq m_i^2 \in L^2(\mathbb{P})$  by Assumptions 11(i, ii, and v). Therefore, the SULLN implies that

$$\sup_{(\gamma, \theta, \xi) \in \Gamma \times \Theta \times \Xi} \left| \frac{1}{n} \sum_{i=1}^n s_i(\xi)\{\tilde{g}_i(\gamma, \xi) - \rho_i(\gamma, \theta)\} \nabla'_{\theta} \rho_i(\gamma, \theta) - \mathbb{E}_{\mathbb{P}}[s_i(\xi)\{\tilde{g}_i(\gamma, \xi) - \rho_i(\gamma, \theta)\} \nabla'_{\theta} \rho_i(\gamma, \theta)] \right| \rightarrow 0$$

a.s.  $-\mathbb{P}$ . Applying the DCT implies that

$$\begin{aligned} \widehat{K}_n &:= \frac{1}{n} \sum_{i=1}^n \int s_i(\hat{\xi}_n)\{\tilde{g}_i(\gamma, \hat{\xi}_n) - \rho_i(\gamma, \tilde{\theta}_n)\} \nabla'_{\theta} \rho_i(\gamma, \tilde{\theta}_n) d\mathbb{Q} \\ &\rightarrow \int E_{\mathbb{P}}[s_i(\xi_*)\{\tilde{g}_i(\gamma, \xi_*) - \rho_i(\gamma, \theta_*)\} \nabla'_{\theta} \rho_i(\gamma, \theta_*)] d\mathbb{Q}(\gamma) =: K \text{ a.s. } -\mathbb{P} \end{aligned}$$

using the fact that  $(\hat{\xi}_n, \tilde{\theta}_n) \rightarrow (\xi_*, \theta_*)$  a.s.  $-\mathbb{P}$ .

The consistency of  $\tilde{B}_n$  for  $B_*$  now follows as a consequence of (i, ii, and iii) in this proof.  $\blacksquare$

**Proof of Lemma 1:** (i) First, Theorem 2(i) implies that  $\sup_{\theta \in \Theta} |q_n(\theta) - q(\theta)| \rightarrow 0$  a.s.  $-\mathbb{P}$ , so that the CFLS estimator  $\hat{\theta}_n^b$  must converge to  $\theta_{\dagger}$  a.s.  $-\mathbb{P}$ , as  $\theta_{\dagger}$  is constrained by the same constraint  $h(\theta) = 0$ . Second,  $\theta_*$  is the global minimizer of  $q(\cdot)$  and also satisfies that  $h(\theta_*) = 0$  under  $\mathbb{H}_o$ . Therefore,  $\theta_* = \theta_{\dagger}$ , as desired.

(ii) Theorem 4(i) implies that  $\sup_{\theta \in \Theta} |\hat{q}_n(\theta) - q(\theta)| \rightarrow 0$  a.s.  $-\mathbb{P}$ , so that the CTSFLS estimator  $\hat{\theta}_n^{\ddagger}$  must converge to  $\theta_{\dagger}$  a.s.  $-\mathbb{P}$ , as  $\theta_{\dagger}$  is constrained by the same constraint  $h(\theta) = 0$ . The remainder of the proof is identical to the proof of Lemma 1(i).  $\blacksquare$

**Proof of Theorem 8:** (i) By virtue of the mean-value theorem, for some  $\bar{\theta}_n$  between  $\hat{\theta}_n$  and  $\theta_*$ ,  $\sqrt{n}\{h(\hat{\theta}_n) - h(\theta_*)\} = D(\bar{\theta}_n)\sqrt{n}(\hat{\theta}_n - \theta_*)$ , so that Theorems 2 and 3 imply that  $\sqrt{n}\{h(\hat{\theta}_n) - h(\theta_*)\} \overset{\Delta}{\rightsquigarrow} \mathcal{N}(0, D_* A^{-1} B A^{-1} D'_*)$ . Further,  $\hat{B}_n \rightarrow B$  a.s.  $-\mathbb{P}$  by Theorem 6, so that  $n\{h(\hat{\theta}_n) - h(\theta_*)\}' \{D_* A^{-1} \hat{B}_n A^{-1} D'_*\}^{-1} \{h(\hat{\theta}_n) - h(\theta_*)\} \overset{\Delta}{\rightsquigarrow} \mathcal{X}^2(r, 0)$ . Therefore, under  $\mathbb{H}_o$ ,  $\mathcal{W}_n^b \overset{\Delta}{\rightsquigarrow} \mathcal{X}^2(r, 0)$ . Meanwhile,  $nh(\hat{\theta}_n)' \{D_* A^{-1} \hat{B}_n A^{-1} D'_*\}^{-1} h(\hat{\theta}_n) = O_{\mathbb{P}}(n)$  but not  $o_{\mathbb{P}}(n)$  because  $h(\hat{\theta}_n) \rightarrow h(\theta_*) \neq 0$  under  $\mathbb{H}_a$ , so that the desired result follows.

(ii) The proofs are almost identical to Theorem 8(i): under  $\mathbb{H}_o$ ,  $\sqrt{n}\{h(\tilde{\theta}_n) - h(\theta_*)\} \overset{\Delta}{\rightsquigarrow} \mathcal{N}(0, D_* A^{-1} B_* A^{-1} D'_*)$  by the mean-value theorem and Theorems 4 and 5; and  $\tilde{B}_n \rightarrow B_*$  a.s.  $-\mathbb{P}$  by Theorem 7, so that  $n\{h(\tilde{\theta}_n) - h(\theta_*)\}' \{D_* A^{-1} \tilde{B}_n A^{-1} D'_*\}^{-1} \{h(\tilde{\theta}_n) - h(\theta_*)\} \overset{\Delta}{\rightsquigarrow} \mathcal{X}^2(r, 0)$  under  $\mathbb{H}_o$ , implying that  $\mathcal{W}_n^{\ddagger} \overset{\Delta}{\rightsquigarrow} \mathcal{X}^2(r, 0)$  under  $\mathbb{H}_o$ . Meanwhile,  $nh(\tilde{\theta}_n)' \{D_* A^{-1} \tilde{B}_n$

$A^{-1}D'_* \}^{-1}h(\tilde{\theta}_n) = O_{\mathbb{P}}(n)$  but not  $o_{\mathbb{P}}(n)$  because  $h(\tilde{\theta}_n) \rightarrow h(\theta_*) \neq 0$  under  $\mathbb{H}_a$ , as desired.  $\blacksquare$

**Proof of Theorem 9:** (i) The CFLS estimator can be obtained by minimizing the Lagrange function:  $\mathcal{L}_n(\theta, \lambda) := q_n(\theta) - \lambda' h(\theta)$ , whose first-order conditions can be given as  $\nabla_{\theta} q_n(\tilde{\theta}_n^b) - \check{\lambda}_n^{b'} \check{D}_n^b \equiv 0$  and  $h(\tilde{\theta}_n^b) \equiv 0$ , where  $(\check{\theta}_n^{b'}, \check{\lambda}_n^{b'})'$  is the solution for the first-order conditions. In addition, for some  $\bar{\theta}_n^b$  between  $\tilde{\theta}_n^b$  and  $\theta_*$ ,

$$\nabla_{\theta} q_n(\tilde{\theta}_n^b) = \nabla_{\theta} q_n(\theta_*) + \nabla_{\theta}^2 q_n(\bar{\theta}_n^b)(\tilde{\theta}_n^b - \theta_*) \text{ and } h(\tilde{\theta}_n^b) = h(\theta_*) + \check{D}_n^b(\tilde{\theta}_n^b - \theta_*) + o_{\mathbb{P}}(1).$$

Plugging these into the first-order conditions, we obtain that

$$\begin{aligned} \sqrt{n}(\tilde{\theta}_n^b - \theta_*) &= -\{\nabla_{\theta}^2 q_n(\bar{\theta}_n^b)\}^{-1}(J - \check{D}_n^{b'} \{\check{E}_n^b\}^{-1} \check{D}_n^b \{\nabla_{\theta}^2 q_n(\bar{\theta}_n^b)\}^{-1}) \sqrt{n} \psi_n \\ &\quad - \{\nabla_{\theta}^2 q_n(\bar{\theta}_n^b)\}^{-1} \check{D}_n^{b'} \{\check{E}_n^b\}^{-1} \sqrt{n} h(\theta_*) + o_{\mathbb{P}}(1) \quad \text{and} \end{aligned} \quad (33)$$

$$\sqrt{n} \check{\lambda}_n^b = \{\check{E}_n^b\}^{-1} \check{D}_n^b \{\nabla_{\theta}^2 q_n(\bar{\theta}_n^b)\}^{-1} \sqrt{n} \psi_n + \{\check{E}_n^b\}^{-1} \sqrt{n} h(\theta_*),$$

where  $\check{E}_n^b := -\check{D}_n^b \{\nabla_{\theta}^2 q_n(\bar{\theta}_n^b)\}^{-1} \check{D}_n^{b'}$ , and  $\psi_n := \nabla_{\theta} q_n(\theta_*)$ .

Given this, applying Theorem 3 implies that  $\sqrt{n} \psi_n \overset{\Delta}{\rightsquigarrow} \mathcal{N}(0, 4B)$ . Furthermore, under  $\mathbb{H}_o$ ,  $\tilde{\theta}_n^b \rightarrow \theta_*$  a.s.- $\mathbb{P}$ , and  $\hat{\theta}_n \rightarrow \theta_*$  a.s.- $\mathbb{P}$ , so that  $\check{E}_n^b \rightarrow -\frac{1}{2} D_* A^{-1} D'_*$  a.s.- $\mathbb{P}$ . Therefore,  $\sqrt{n} \check{\lambda}_n^b \overset{\Delta}{\rightsquigarrow} N[0, 4(D_* A^{-1} D'_*)^{-1} D_* A^{-1} B A^{-1} D'_* (D_* A^{-1} D'_*)^{-1}]$ , implying that  $\frac{n}{4} \check{\lambda}_n^{b'} D_* A^{-1} D'_* (D_* A^{-1} B A^{-1} D'_*)^{-1} D_* A^{-1} D'_* \check{\lambda}_n^b \overset{\Delta}{\rightsquigarrow} \mathcal{X}^2(r, 0)$ . Given this result, we also obtain that  $\mathcal{LM}_n^b \overset{\Delta}{\rightsquigarrow} \mathcal{X}^2(r, 0)$  from Theorem 6,  $\nabla_{\theta} q_n(\tilde{\theta}_n^b) \equiv \check{\lambda}_n^{b'} \check{D}_n^b$ , and the fact that  $\tilde{\theta}_n^b \rightarrow \theta_*$  a.s.- $\mathbb{P}$  under  $\mathbb{H}_o$ . Meanwhile,  $\sqrt{n} h(\theta_*) = O_{\mathbb{P}}(\sqrt{n})$  though not  $o_{\mathbb{P}}(\sqrt{n})$  under  $\mathbb{H}_a$ ;  $\check{D}_n^b \rightarrow D_{\dagger} := D(\theta_{\dagger})$  a.s.- $\mathbb{P}$ ; and  $\nabla_{\theta}^2 q_n(\bar{\theta}_n^b) \rightarrow \nabla_{\theta}^2 q(\theta_{\S})$  a.s.- $\mathbb{P}$  for some  $\theta_{\S}$  such that  $\nabla_{\theta} q(\theta_{\dagger}) = \nabla_{\theta} q(\theta_*) + \nabla_{\theta}^2 q(\theta_{\S})(\theta_{\dagger} - \theta_*)$ , so that  $\sqrt{n} \check{\lambda}_n^b = O_{\mathbb{P}}(\sqrt{n})$  though not  $o_{\mathbb{P}}(\sqrt{n})$ . Therefore, the desired result follows.

(ii) The proofs are almost identical to those of Theorem 9(i). The TSCFLS estimator can be obtained by minimizing the Lagrange function:  $\hat{\mathcal{L}}_n(\theta, \lambda) := \hat{q}_n(\theta) - \lambda' h(\theta)$ , whose first-order conditions are given as  $\nabla_{\theta} \hat{q}_n(\tilde{\theta}_n^{\#}) - \check{\lambda}_n^{\#'} \check{D}_n^{\#} \equiv 0$  and  $h(\tilde{\theta}_n^{\#}) \equiv 0$ , where  $(\check{\theta}_n^{\#'}, \check{\lambda}_n^{\#'})'$  is the solution for the first-order conditions. In addition, for some  $\bar{\theta}_n^{\#}$  between  $\tilde{\theta}_n^{\#}$  and  $\theta_*$ ,

$$\nabla_{\theta} \hat{q}_n(\tilde{\theta}_n^{\#}) = \nabla_{\theta} \hat{q}_n(\theta_*) + \nabla_{\theta}^2 \hat{q}_n(\bar{\theta}_n^{\#})(\tilde{\theta}_n^{\#} - \theta_*) \text{ and } h(\tilde{\theta}_n^{\#}) = h(\theta_*) + \check{D}_n^{\#}(\tilde{\theta}_n^{\#} - \theta_*) + o_{\mathbb{P}}(1).$$

By plugging these into the first-order conditions, we obtain that

$$\begin{aligned} \sqrt{n}(\tilde{\theta}_n^{\#} - \theta_*) &= -\{\nabla_{\theta}^2 \hat{q}_n(\bar{\theta}_n^{\#})\}^{-1}(J - \check{D}_n^{\#'} \{\check{E}_n^{\#}\}^{-1} \check{D}_n^{\#} \{\nabla_{\theta}^2 \hat{q}_n(\bar{\theta}_n^{\#})\}^{-1}) \sqrt{n} \hat{\psi}_n \\ &\quad - \{\nabla_{\theta}^2 \hat{q}_n(\bar{\theta}_n^{\#})\}^{-1} \check{D}_n^{\#'} \{\check{E}_n^{\#}\}^{-1} \sqrt{n} h(\theta_*) + o_{\mathbb{P}}(1) \quad \text{and} \end{aligned} \quad (34)$$

$$\sqrt{n} \check{\lambda}_n^{\#} = \{\check{E}_n^{\#}\}^{-1} \check{D}_n^{\#} \{\nabla_{\theta}^2 \hat{q}_n(\bar{\theta}_n^{\#})\}^{-1} \sqrt{n} \hat{\psi}_n + \{\check{E}_n^{\#}\}^{-1} \sqrt{n} h(\theta_*), \quad (35)$$

where  $\check{E}_n^{\#} := -\check{D}_n^{\#} \{\nabla_{\theta}^2 \hat{q}_n(\bar{\theta}_n^{\#})\}^{-1} \check{D}_n^{\#'}$ , and  $\hat{\psi}_n := \nabla_{\theta} \hat{q}_n(\theta_*)$ .

Given this,  $\sqrt{n} \hat{\psi}_n \overset{\Delta}{\rightsquigarrow} \mathcal{N}(0, 4B_*)$  by applying (27). Further,  $\tilde{\theta}_n^{\#} \rightarrow \theta_*$  a.s.- $\mathbb{P}$  under  $\mathbb{H}_o$ , and  $\tilde{\theta}_n \rightarrow \theta_*$  a.s.- $\mathbb{P}$ , so that  $\check{E}_n^{\#} \rightarrow -\frac{1}{2} D_* A^{-1} D'_*$  a.s.- $\mathbb{P}$  and  $\sqrt{n} \check{\lambda}_n^{\#} \overset{\Delta}{\rightsquigarrow} N[0, 4(D_* A^{-1} D'_*)^{-1} D_* A^{-1} B_* A^{-1} D'_* (D_* A^{-1} D'_*)^{-1}]$ , implying that  $\frac{n}{4} \check{\lambda}_n^{\#'} D_* A^{-1} D'_* (D_* A^{-1} B_* A^{-1} D'_*)^{-1} D_* A^{-1} D'_* \check{\lambda}_n^{\#} \overset{\Delta}{\rightsquigarrow} \mathcal{X}^2(r, 0)$ . Therefore,  $\mathcal{LM}_n^{\#} \overset{\Delta}{\rightsquigarrow} \mathcal{X}^2(r, 0)$  by Theorem 7,  $\nabla_{\theta} \hat{q}_n(\tilde{\theta}_n^{\#}) \equiv \check{\lambda}_n^{\#'} \check{D}_n^{\#}$ , and the fact that  $\tilde{\theta}_n^{\#} \rightarrow \theta_*$  a.s.- $\mathbb{P}$  under  $\mathbb{H}_o$ . Meanwhile,  $\sqrt{n} h(\theta_*) = O_{\mathbb{P}}(\sqrt{n})$  though not  $o_{\mathbb{P}}(\sqrt{n})$  under  $\mathbb{H}_a$ ;  $\check{D}_n^{\#} \rightarrow D_{\dagger} := D(\theta_{\dagger})$  a.s.- $\mathbb{P}$ ; and  $\nabla_{\theta}^2 \hat{q}_n(\bar{\theta}_n^{\#}) \rightarrow \nabla_{\theta}^2 q(\theta_{\S})$  a.s.- $\mathbb{P}$  for some  $\theta_{\S}$  such that  $\nabla_{\theta} q(\theta_{\dagger}) = \nabla_{\theta} q(\theta_*) + \nabla_{\theta}^2 q(\theta_{\S})(\theta_{\dagger} - \theta_*)$ , so that  $\sqrt{n} \check{\lambda}_n^{\#} = O_{\mathbb{P}}(\sqrt{n})$  though



not  $o_{\mathbb{P}}(\sqrt{n})$ . The desired result then follows.  $\blacksquare$

**Proof of Theorem 10:** (i) By the mean-value theorem and the first-order condition for  $\widehat{\theta}_n$ , note that for some  $\widetilde{\theta}_n^b$  between  $\widehat{\theta}_n$  and  $\check{\theta}_n^b$ ,

$$q_n(\check{\theta}_n^b) = q_n(\widehat{\theta}_n) + \frac{1}{2}(\check{\theta}_n^b - \widehat{\theta}_n)' \{ \nabla_{\theta}^2 q_n(\widetilde{\theta}_n^b) \} (\check{\theta}_n^b - \widehat{\theta}_n). \quad (36)$$

Furthermore, it follows that

$$\sqrt{n}(\check{\theta}_n^b - \widehat{\theta}_n) = \{ \nabla_{\theta}^2 q_n(\widetilde{\theta}_n^b) \}^{-1} \check{D}_n^{b'} \{ \check{E}_n^b \}^{-1} \sqrt{n} [ \check{D}_n^b \{ \nabla_{\theta}^2 q_n(\widetilde{\theta}_n^b) \}^{-1} \psi_n - h(\theta_*) ] + o_{\mathbb{P}}(1) \quad (37)$$

from (33) and (22). As given in the proof of Theorem 9(i),  $\check{\theta}_n^b \rightarrow \theta_*$  a.s.- $\mathbb{P}$ ,  $\widetilde{\theta}_n^b \rightarrow \theta_*$  a.s.- $\mathbb{P}$  under  $\mathbb{H}_o$ , and  $\widehat{\theta}_n \rightarrow \theta_*$  a.s.- $\mathbb{P}$ , so that  $\check{E}_n^b \rightarrow -\frac{1}{2}D_*A^{-1}D_*'$  a.s.- $\mathbb{P}$  and  $\nabla_{\theta}^2 q_n(\widetilde{\theta}_n^b) \rightarrow 2A$  a.s.- $\mathbb{P}$ . In addition,  $\sqrt{n}\psi_n \overset{\Delta}{\sim} \mathcal{N}(0, 4B)$  as given in the proof of Theorem 3. Therefore, if we employ all these results in (36),  $n\{q_n(\check{\theta}_n^b) - q_n(\widehat{\theta}_n)\} \Rightarrow W'(D_*A^{-1}D_*')^{-1}W$  under  $\mathbb{H}_o$ . Meanwhile,  $\sqrt{n}(\check{\theta}_n^b - \widehat{\theta}_n)$  is not bounded in probability under  $\mathbb{H}_a$  mainly because  $\sqrt{n}h(\theta_*) = O(\sqrt{n})$  though not  $o(\sqrt{n})$  in (37);  $\check{D}_n^b \rightarrow D_{\dagger} := D(\theta_{\dagger})$  a.s.- $\mathbb{P}$ ; and  $\nabla_{\theta}^2 q_n(\widetilde{\theta}_n^b) \rightarrow \nabla_{\theta}^2 q(\theta_{\S})$  a.s.- $\mathbb{P}$  for some  $\theta_{\S}$  such that  $\nabla_{\theta} q(\check{\theta}_n^b) = \nabla_{\theta} q(\theta_*) + \nabla_{\theta}^2 q(\theta_{\S})(\check{\theta}_n^b - \theta_*)$ . The desired result then follows.

(ii) The proofs follow those of Theorem 10(i). By the mean-value theorem and the first-order condition for  $\widetilde{\theta}_n$ , for some  $\check{\theta}_n^{\#}$  between  $\widetilde{\theta}_n$  and  $\check{\theta}_n^{\#}$ ,

$$q_n(\check{\theta}_n^{\#}) = q_n(\widetilde{\theta}_n) + \frac{1}{2}(\check{\theta}_n^{\#} - \widetilde{\theta}_n)' \{ \nabla_{\theta}^2 q_n(\widetilde{\theta}_n^{\#}) \} (\check{\theta}_n^{\#} - \widetilde{\theta}_n). \quad (38)$$

Furthermore, it follows that

$$\sqrt{n}(\check{\theta}_n^{\#} - \widetilde{\theta}_n) = \{ \nabla_{\theta}^2 q_n(\widetilde{\theta}_n^{\#}) \}^{-1} \check{D}_n^{\#'} \{ \check{E}_n^{\#} \}^{-1} \sqrt{n} [ \check{D}_n^{\#} \{ \nabla_{\theta}^2 q_n(\widetilde{\theta}_n^{\#}) \}^{-1} \widehat{\psi}_n - h(\theta_*) ] + o_{\mathbb{P}}(1), \quad (39)$$

by (34) and (25). As given in the proof of Theorem 9(ii), we have  $\check{\theta}_n^{\#} \rightarrow \theta_*$  a.s.- $\mathbb{P}$ ,  $\widetilde{\theta}_n^{\#} \rightarrow \theta_*$  a.s.- $\mathbb{P}$  under  $\mathbb{H}_o$ , and  $\widetilde{\theta}_n \rightarrow \theta_*$  a.s.- $\mathbb{P}$ , so that  $\check{E}_n^{\#} \rightarrow -\frac{1}{2}D_*A^{-1}D_*'$  a.s.- $\mathbb{P}$  and  $\nabla_{\theta}^2 q_n(\widetilde{\theta}_n^{\#}) \rightarrow 2A$  a.s.- $\mathbb{P}$ . Further,  $\sqrt{n}\widehat{\psi}_n \overset{\Delta}{\sim} \mathcal{N}(0, 4B_*)$  as given in the proof of Theorem 5. Thus, if we use these results in (38),  $n\{q_n(\check{\theta}_n^{\#}) - q_n(\widetilde{\theta}_n)\} \Rightarrow W_*'(D_*A^{-1}D_*')^{-1}W_*$  under  $\mathbb{H}_o$ . Meanwhile,  $\sqrt{n}(\check{\theta}_n^{\#} - \widetilde{\theta}_n)$  is not bounded in probability under  $\mathbb{H}_a$  mainly because  $\sqrt{n}h(\theta_*) = O(\sqrt{n})$  though not  $o(\sqrt{n})$  in (39);  $\check{D}_n^{\#} \rightarrow D_{\dagger} := D(\theta_{\dagger})$  a.s.- $\mathbb{P}$ ; and  $\nabla_{\theta}^2 q_n(\widetilde{\theta}_n^{\#}) \rightarrow \nabla_{\theta}^2 q(\theta_{\S})$  a.s.- $\mathbb{P}$  for some  $\theta_{\S}$  such that  $\nabla_{\theta} q(\check{\theta}_n^{\#}) = \nabla_{\theta} q(\theta_*) + \nabla_{\theta}^2 q(\theta_{\S})(\check{\theta}_n^{\#} - \theta_*)$ . The desired result follows.  $\blacksquare$

## C Estimating the Population Mean Function

This section explores estimation of the population mean function of  $g_i(\cdot)$ . For each  $\gamma$  the mean quantity

$$\mu(\gamma) := \int g(\gamma) d\mathbb{P}(g(\gamma))$$

is no longer a function of  $x$  and for each  $\gamma$  we may denote  $\mu(\gamma)$  as  $\mathbb{E}[g_i(\gamma)]$ . Estimation and inference on  $\mu(\cdot)$  cannot be made by  $\mathcal{M}$  due to the presence of  $x_i$  in  $\rho_i(\cdot, \theta)$ . We, therefore, suppose another model without  $x_i$  as follows:

$$\mathcal{M}_0 := \{ \rho(\cdot, \theta) : \Gamma \mapsto \mathbb{R} | \theta \in \Theta \in \mathbb{R}^d \}$$

and estimate  $\mu(\cdot)$  by FLS estimation, i.e.,

$$\ddot{\theta}_n := \arg \min_{\theta \in \Theta} \ddot{q}_n(\theta), \quad \text{where} \quad \ddot{q}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \int \{\tilde{g}_i(\gamma, \hat{\xi}_n) - \rho(\gamma, \theta)\}^2 d\mathbb{Q}(\gamma),$$

which again is designed to consistently estimate

$$\ddot{\theta} := \arg \min_{\theta \in \Theta} \ddot{q}(\theta), \quad \text{where} \quad \ddot{q}(\theta) := \int \int \{g(\gamma) - \rho(\gamma, \theta)\}^2 d\mathbb{P}(g(\gamma)) d\mathbb{Q}(\gamma).$$

Asymptotic analysis of this FLS estimator is a special case of the analysis in Section 3, so that consistency and asymptotic normality of  $\ddot{\theta}_n$  can be achieved under milder conditions than those of Section 3. We first collect the conditions together as follows.

**Assumption 13.** (i) For each  $\theta \in \Theta$ ,  $\rho(\cdot, \theta) : \Gamma \mapsto \mathbb{R}$  is measurable  $-G$ ;

(ii) for each  $\gamma \in \Gamma$ ,  $\rho(\gamma, \cdot) : \Theta \mapsto \mathbb{R}$  is in  $\mathcal{C}^{(2)}(\Theta)$ ;

(iii)  $\Theta$  is a compact and convex set in  $\mathbb{R}^d$  ( $d \in \mathbb{N}$ );

(iv)  $\ddot{\theta}_*$  is unique and lies in the interior of  $\Theta$ ;

(v)  $A_0$  is positive definite, where  $A_0 := \int \nabla_{\theta} \rho(\gamma, \ddot{\theta}_*) \nabla'_{\theta} \rho(\gamma, \ddot{\theta}_*) d\mathbb{Q}(\gamma) - \int \{\mu(\gamma) - \rho(\gamma, \theta_*)\} \nabla_{\theta}^2 \rho(\gamma, \ddot{\theta}_*) d\mathbb{Q}(\gamma)$ ;

(vi) for some  $m_i \in L^2(\mathbb{P})$ ,  $\sup_{(\gamma, \xi) \in \Gamma \times \Xi} |\tilde{g}_i(\gamma, \xi)| \leq m_i$  a.s.  $-\mathbb{P}$  and  $\sup_j \sup_{(\gamma, \xi) \in \Gamma \times \Xi} |(\partial/\partial \xi_j) \tilde{g}_i(\gamma, \xi)| \leq m_i$  a.s.  $-\mathbb{P}$ ;

(vii)  $\sup_{\theta \in \Theta} |\rho(\cdot, \theta)| \in L^2(\mathbb{Q})$  and for each  $j$  and  $j' = 1, 2, \dots, d$ ,  $\sup_{\theta \in \Theta} |(\partial/\partial \theta_j) \rho(\cdot, \theta)| \in L^2(\mathbb{Q})$  and  $\sup_{\theta \in \Theta} |(\partial^2/\partial \theta_j \partial \theta_{j'}) \rho(\cdot, \theta)| \in L^2(\mathbb{Q})$ ;

(viii)  $C_0$  is positive definite, where

$$C_0 := \begin{bmatrix} J & K'_0 \\ K_0 & B_0 \end{bmatrix},$$

$K_0 := \int \mathbb{E}_{\mathbb{P}}[s_i(\xi_*) \{\tilde{g}_i(\gamma, \xi_*) - \rho(\gamma, \ddot{\theta}_*)\}] \nabla'_{\theta} \rho(\gamma, \ddot{\theta}_*) d\mathbb{Q}(\gamma)$ ,  $B_0 := \int \int \nabla_{\theta} \rho(\gamma, \ddot{\theta}_*) \kappa(\gamma, \tilde{\gamma}) \nabla'_{\theta} \rho(\tilde{\gamma}, \ddot{\theta}_*) d\mathbb{Q}(\gamma) d\mathbb{Q}(\tilde{\gamma})$ , and  $\kappa(\gamma, \tilde{\gamma}) := \int (g(\gamma) - \rho(\gamma, \ddot{\theta}_*)) (g(\tilde{\gamma}) - \rho(\tilde{\gamma}, \ddot{\theta}_*)) d\mathbb{P}(g(\gamma), g(\tilde{\gamma}))$ ; and

(ix)  $B_{\dagger}$  is positive definite, where  $B_{\dagger} := B_0 - M_0 H^{-1} K_0 - K'_0{}^{-1} M_0 + M_0 H^{-1} J H^{-1} M'_0$  and  $M_0 := \int \nabla_{\theta} \rho(\gamma, \ddot{\theta}_*) \mathbb{E}_{\mathbb{P}}[\nabla'_{\xi} \tilde{g}_i(\gamma, \xi_*)] d\mathbb{Q}(\gamma)$ .  $\square$

There is a correspondence between Assumption 13 and the earlier conditions. Assumptions 13(i, ii, iii, and iv) correspond to Assumption 2 and, due to the absence of  $x_i$  in the model, the conditions in Assumption 2 are appropriately modified. Assumption 13(iv) also corresponds to Assumption 4. Note that the integrands of  $A_0$  are non-stochastic, whereas  $A$  has stochastic integrands. This difference again stems from the absence of  $x_i$  from the model  $\mathcal{M}_0$ . Assumptions 13(vi and vii) correspond to Assumption 7. Assumption 13(vi) is the same as Assumptions 7(i and iii), but Assumption 13(vii) is milder than Assumption 7(ii, iv, and v). We do not need the stochastic bound conditions for the model and its derivatives due to the absence of  $x_i$ . Finally, note that Assumptions 13(viii and ix) correspond to Assumption 9.

The following corollary gives the limit behavior of  $\ddot{q}_n(\cdot)$  and  $\ddot{\theta}_n$  using Assumption 13 in addition to the conditions for  $\hat{\xi}_n$ .

**Corollary (A).** Given Assumptions 6, 8, and 13,

(i)  $\sup_{\theta \in \Theta} |\ddot{q}_n(\theta) - \ddot{q}(\theta)| \rightarrow 0$  a.s.  $-\mathbb{P}$ ;

(ii)  $\ddot{\theta}_n \rightarrow \ddot{\theta}_*$  a.s.  $-\mathbb{P}$ ;

(iii)  $\sqrt{n}(\ddot{\theta}_n - \ddot{\theta}_*) \overset{A}{\rightsquigarrow} \mathcal{N}(A_0^{-1} B_{\dagger} A_0^{-1})$ ; and

(iv) if  $\xi_*$  is known,  $\sqrt{n}(\ddot{\theta}_n - \ddot{\theta}_*) \overset{A}{\rightsquigarrow} \mathcal{N}(A_0^{-1}B_0A_0^{-1})$ . □

In view of the parallel structure of  $\ddot{\theta}_n$  to  $\tilde{\theta}_n$ , Corollary (C) can be established by repeating the arguments given in earlier sections. The proof is omitted.

## D Derivation of Test Statistic Basis in Section 5.4

The form of  $g_i$  is derived by following Davies (1977). Specifically, it is obtained by imposing the hypothesis  $\pi_* = 0$  to the standardized score obtained with respect to  $\pi$ . That is,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n g_i(U_i, V_i) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial}{\partial \pi} \log \{(1 - \pi)c_1(U_i, V_i; \gamma_1) + \pi c_2(U_i, V_i; \gamma_2)\} / W(\gamma_1, \gamma_2, \pi) \Big|_{\pi=0},$$

where

$$W(\gamma_1, \gamma_2, \pi) := \text{var} \left[ \frac{\partial}{\partial \pi} \log \{(1 - \pi)c_1(U_i, V_i; \gamma_1) + \pi c_2(U_i, V_i; \gamma_2)\} \right]^{1/2}.$$

In our setup note that

$$\frac{\partial}{\partial \pi} \log \{(1 - \pi)c_1(U_i, V_i; \gamma_1) + \pi c_2(U_i, V_i; \gamma_2)\} \Big|_{\pi=0} = \frac{c_2(U_i, V_i; \gamma_2) - c_1(U_i, V_i; \gamma_1)}{c_1(U_i, V_i; \gamma_1)}$$

and

$$W(\gamma_1, \gamma_2, 1) = \sqrt{\int_0^1 \int_0^1 \frac{c_2^2(u, v; \gamma_2)}{c_1(u, v; \gamma_1)} dudv} - 1,$$

leading to the desired expression for  $g_i$  in Section 5.4.

## References

- AI, C. AND X. CHEN (2003): “Efficient Estimation of Models with Conditional Moment Restrictions Containing Unknown Functions,” *Econometrica*, 71, 1795–1843.
- AKHARIF, A., M. FIHRI, M. HALLIN, AND A. MELLOUK (2020): “Optimal Pseudo-Gaussian and Rank-Based Random Coefficient Detection in Multiple Regression,” *Electronic Journal of Statistics*, 14, 4207 – 4243.
- AMEMIYA, T. (1979): “The Estimation of a Simultaneous-Equation Tobit Model,” *International Economic Review*, 20, 169–181.
- ANDRES, C., A. BETZER, M. GOERGEN, AND L. RENNEBOOG (2009): “Dividend Policy of German Firms: A Panel Data Analysis of Partial Adjustment Models,” *Journal of Empirical Finance*, 16, 175–187.
- ANDREWS, D. W. K. (1992): “Generic Uniform Convergence,” *Econometric Theory*, 8, 241–257.
- (2001): “Testing When a Parameter is on the Boundary of the Maintained Hypothesis,” *Econometrica*, 69, 683–734.
- AUTOR, D. H., L. F. KATZ, AND A. B. KRUEGER (1998): “Computing Inequality: Have Computers Changed the Labor Market?” *The Quarterly Journal of Economics*, 113, 1169–1213.

- BAEK, Y. I., J. S. CHO, AND P. C. B. PHILLIPS (2015): “Testing Linearity Using Power Transforms of Regressors,” *Journal of Econometrics*, 187, 376–384.
- BARTH, E., J. DAVIS, AND R. B. FREEMAN (2018): “Augmenting the Human Capital Earnings Equation with Measures of Where People Work,” *Journal of Labor Economics*, 36, 71–97.
- BEARE, B. K., J. SEO, AND W.-K. SEO (2017): “Cointegrated Linear Processes in Hilbert Space,” *Journal of Time Series Analysis*, 38, 1010–1027.
- BHULLER, M., M. MOGSTAD, AND K. G. SALVANES (2017): “Life-Cycle Earnings, Education Premiums, and Internal Rates of Return,” *Journal of Labor Economics*, 35, 993–1030.
- BIERENS, H. J. (1990): “A Consistent Conditional Moment Test of Functional Form,” *Econometrica*, 58, 1443–1458.
- BOSQ, D. (2000): *Linear Processes in Function Spaces: Theory and Applications*, New York: Springer.
- BRAV, A., J. R. GRAHAM, C. R. HARVEY, AND R. MICHAELY (2005): “Payout Policy in the 21st Century,” *Journal of Financial Economics*, 77, 483–527.
- BREUSCH, T. S. AND A. R. PAGAN (1979): “A Simple Test for Heteroscedasticity and Random Coefficient Variation,” *Econometrica*, 47, 1287–1294.
- BUGNI, F. A., P. HALL, J. L. HOROWITZ, AND G. R. NEUMANN (2009): “Goodness-of-Fit Tests for Functional Data,” *The Econometrics Journal*, 12, 1–18.
- CAI, T. T. AND P. HALL (2006): “Prediction in Functional Linear Regression,” *The Annals of Statistics*, 34, 2159–2179.
- CAO, G., L. YANG, AND D. TODEM (2012): “Simultaneous Inference for the Mean Function Based on Dense Functional Data,” *Journal of Nonparametric Statistics*, 24, 359–377.
- CARDOT, H., C. CRAMBES, A. KNEIP, AND P. SARDA (2007): “Smoothing Splines Estimators in Functional Linear Regression with Errors-in-Variables,” *Computational Statistics & Data Analysis*, 51, 4832–4848.
- CHANG, Y., B. HU, AND J. Y. PARK (2019): “Econometric Analysis of Functional Dynamics in the Presence of Persistence,” Tech. rep., Indiana University.
- CHANG, Y., C. S. KIM, AND J. Y. PARK (2016): “Nonstationarity in Time Series of State Densities,” *Journal of Econometrics*, 192, 152–167.
- CHEN, H., J. CHEN, AND J. D. KALBFLEISCH (2001): “A Modified Likelihood Ratio Test for Homogeneity in Finite Mixture Models,” *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 63, 19–29.
- CHEN, J. AND P. LI (2009): “Hypothesis Test for Normal Mixture Models: The EM Approach,” *The Annals of Statistics*, 37, 2523 – 2542.
- CHEN, X. (2007): “Large Sample Sieve Estimation of Semi-Nonparametric Models,” in *Handbook of Econometrics*, ed. by J. J. Heckman and E. E. Leamer, Amsterdam: North Holland. Chapter 76, pp. 5549–5632, vol. 6.

- CHEN, X. AND Y. FAN (2006): “Estimation and Model Selection of Semiparametric Copula-Based Multivariate Dynamic Models under Copula Misspecification,” *Journal of Econometrics*, 135, 125–154.
- CHERNOFF, H. AND E. LANDER (1995): “Asymptotic Distribution of the Likelihood Ratio Test That a Mixture of Two Binomials is a Single Binomial,” *Journal of Statistical Planning and Inference*, 43, 19–40.
- CHILÈS, J.-P. AND P. DELFINER (1999): *Geostatistics: Modelling Spatial Uncertainty*, New York: John Wiley & Sons.
- CHO, J. S., M.-H. PARK, AND P. C. B. PHILLIPS (2018): “Practical Kolmogorov–Smirnov Testing by Minimum Distance Applied to Measure Top Income Shares in Korea,” *Journal of Business & Economic Statistics*, 36, 523–537.
- CHO, J. S. AND P. C. B. PHILLIPS (2018a): “Pythagorean Generalization of Testing the Equality of Two Symmetric Positive Definite Matrices,” *Journal of Econometrics*, 202, 45–56.
- (2018b): “Sequentially Testing Polynomial Model Hypotheses Using Power Transforms of Regressors,” *Journal of Applied Econometrics*, 33, 141–159.
- CHO, J. S. AND H. WHITE (2007): “Testing for Regime Switching,” *Econometrica*, 75, 1671–1720.
- (2010): “Testing for Unobserved Heterogeneity in Exponential and Weibull Duration Models,” *Journal of Econometrics*, 157, 458–480.
- (2014): “Testing the Equality of Two Positive-Definite Matrices with Application to Information Matrix Testing,” in *Advances in Econometrics: Essays in Honor of Peter C. B. Phillips*, ed. by Y. Chang, T. B. Fomby, and J. Y. Park, West Yorkshire, UK: Emerald Group Publishing Limited. pp. 491–556, vol. 33.
- CRAMBES, C., A. GANNOUN, AND Y. HENCHIRI (2013): “Support Vector Machine Quantile Regression Approach for Functional Data: Simulation and Application Studies,” *Journal of Multivariate Analysis*, 121, 50–68.
- DAVIES, R. B. (1977): “Hypothesis Testing When a Nuisance Parameter is Present Only Under the Alternative,” *Biometrika*, 64, 247–254.
- (1987): “Hypothesis Testing when a Nuisance Parameter is Present Only Under the Alternatives,” *Biometrika*, 74, 33–43.
- DIAS, A. AND P. EMBRECHTS (2004): “Dynamic Copula Models for Multivariate High-Frequency Data in Finance,” Tech. Rep. 04-01, Warwick Business School, Finance Group.
- DIKS, C., V. PANCHENKO, AND D. VAN DIJK (2010): “Out-of-Sample Comparison of Copula Specifications in Multivariate Density Forecasts,” *Journal of Economic Dynamics and Control*, 34, 1596–1609.
- EFRON, B. (1979): “Bootstrap Methods: Another Look at the Jackknife,” *The Annals of Statistics*, 7, 1–26.
- ENGLE, R. AND M. WATSON (1985): “The Kalman Filter: Applications to Forecasting and Rational Expectations Models,” in *Advances in Econometrics: Fifth World Congress*, ed. by T. Bewley, New York: Cambridge University Press, vol. 1, 245–293.
- EVERITT, B. AND D. HAND (1981): *Finite Mixture Distributions*, Netherlands: Springer.
- FAMA, E. F. AND H. BABIAK (1968): “Dividend Policy: An Empirical Analysis,” *Journal of the American Statistical Association*, 63, 1132–1161.

- FERMANIAN, J.-D., D. RADULOVIĆ, AND M. WEGKAMP (2004): “Weak Convergence of Empirical Copula Processes,” *Bernoulli*, 10, 847–860.
- FERRATY, F. AND P. VIEU (2006): *Nonparametric Functional Data Analysis: Theory and Practice*, New York: Springer-Verlag.
- FISHER, R. (1932): *Statistical Methods for Research Workers*, Edinburgh and London: Oliver and Boyd.
- FRANCHI, M. AND P. PARUOLO (2020): “Cointegration in Functional Autoregressive Process,” *Econometric Theory*, 36, 803–839.
- FU, Y., J. CHEN, AND P. LI (2008): “Modified Likelihood Ratio Test for Homogeneity in a Mixture of von Mises Distributions,” *Journal of Statistical Planning and Inference*, 138, 667–681.
- GARRETT, I. AND R. PRIESTLEY (2000): “Dividend Behavior and Dividend Signaling,” *The Journal of Financial and Quantitative Analysis*, 35, 173–189.
- GRENANDER, U. (1981): *Abstract Inference*, New York: John Wiley & Sons.
- GUVENEN, F. (2007): “Learning Your Earning: Are Labor Income Shocks Really Very Persistent?” *American Economic Review*, 97, 687–712.
- (2009): “An Empirical Investigation of Labor Income Processes,” *Review of Economic Dynamics*, 12, 58–79.
- HALL, P. AND J. L. HOROWITZ (2007): “Methodology and Convergence Rates for Functional Linear Regression,” *The Annals of Statistics*, 35, 70 – 91.
- HANSEN, B. E. (1996): “Inference When a Nuisance Parameter Is Not Identified Under the Null Hypothesis,” *Econometrica*, 64, 413–430.
- HAPP, C. AND S. GREVEN (2018): “Multivariate Functional Principal Component Analysis for Data Observed on Different (Dimensional) Domains,” *Journal of the American Statistical Association*, 113, 649–659.
- HECKMAN, J. J., L. J. LOCHNER, AND P. E. TODD (2006): “Earnings Functions, Rates of Return and Treatment Effects: The Mincer Equation and Beyond,” in *Handbook of the Economics of Education*, ed. by E. Hanushek and F. Welch, Amsterdam: North Holland, Chapter 7, pp. 307–458, vol. 1.
- HORVATH, L. AND P. KOKOSZKA (2012): *Inference for Functional Data with Applications*, vol. 200, New York: Springer-Verlag.
- HÖRMANN, S., P. KOKOSZKA, AND G. NISOL (2018): “Testing for Periodicity in Functional Time Series,” *The Annals of Statistics*, 46, 2960 – 2984.
- HSIAO, C. (1974): “Statistical Inference for a Model with Both Random Cross-Sectional and Time Effects,” *International Economic Review*, 15, 12–30.
- HU, L. (2006): “Dependence Patterns across Financial Markets: A Mixed Copula Approach,” *Applied Financial Economics*, 16, 717–729.

- HUIZINGA, F. (1990): "An Overlapping Generations Model of Wage Determination," *The Scandinavian Journal of Economics*, 92, 81–98.
- JAMES, G. M., J. WANG, AND J. ZHU (2009): "Functional Linear Regression That's Interpretable," *The Annals of Statistics*, 37, 2083 – 2108.
- JOE, H. (2001): *Multivariate Models and Dependence Concepts*, Monographs on Statistics and Applied Probability 73. London: Taylor & Francis.
- (2014): *Dependence Modeling with Copulas*, New York: Chapman & Hall.
- JOE, H. AND J. XU (1996): "The Estimation Method of Inference Functions for Margins for Multivariate Models," Technical Report 166, Department of Statistics, University of British Columbia.
- KATZ, L. F. AND K. M. MURPHY (1992): "Changes in Relative Wages, 1963-1987: Supply and Demand Factors," *The Quarterly Journal of Economics*, 107, 35–78.
- KIM, J. AND J. Y. PARK (2017): "Asymptotics for Recurrent Diffusions with Application to High Frequency Regression," *Journal of Econometrics*, 196, 37–54.
- KOSMIDIS, I. AND D. KARLIS (2016): "Model-Based Clustering Using Copulas with Applications," *Statistics and Computing*, 26, 1079–1099.
- KUTOYANTS, Y. (1984): *Parameter Estimation for Stochastic Processes*, Berlin: Heldermann Verlag.
- LAI, Y., C. W. S. CHEN, AND R. GERLACH (2009): "Optimal Dynamic Hedging via Copula-Threshold-GARCH Models," *Mathematics and Computers in Simulation*, 79, 2609–2624.
- LANCASTER, H. O. (1961): "The Combination of Probabilities: An Application of Orthonormal Functions," *Australian Journal of Statistics*, 3, 20–33.
- LEE, L.-F., G. S. MADDALA, AND R. P. TROST (1980): "Asymptotic Covariance Matrices of Two-Stage Probit and Two-Stage Tobit Methods for Simultaneous Equations Models with Selectivity," *Econometrica*, 48, 491–503.
- LEMIEUX, T. (2006): "The "Mincer Equation" Thirty Years After *Schooling, Experience, and Earnings*," in *Jacob Mincer A Pioneer of Modern Labor Economics*, ed. by S. Grossbard, Boston: Springer, 127–145.
- LI, D., P. M. ROBINSON, AND H. L. SHANG (2020): "Long-Range Dependent Curve Time Series," *Journal of the American Statistical Association*, 115, 957–971.
- LIANG, K.-Y. AND P. J. RATHOUZ (1999): "Hypothesis Testing Under Mixture Models: Application to Genetic Linkage Analysis," *Biometrics*, 55, 65–74.
- LIGHT, A. AND M. URETA (1995): "Early-Career Work Experience and Gender Wage Differentials," *Journal of Labor Economics*, 13, 121–154.
- LILLARD, L. A. AND Y. WEISS (1979): "Components of Variation in Panel Earnings Data: American Scientists 1960-70," *Econometrica*, 47, 437–454.

- LINTNER, J. (1956): "Distribution of Incomes of Corporations Among Dividends, Retained Earnings, and Taxes," *American Economic Review*, 46, 97–113.
- LOAIZA-MAYA, R., M. S. SMITH, AND W. MANEESOONTHORN (2018): "Time Series Copulas for Heteroskedastic Data," *Journal of Applied Econometrics*, 33, 332–354.
- MAGNAC, T., N. PISTOLESI, AND S. ROUX (2018): "Post-Schooling Human Capital Investments and the Life Cycle of Earnings," *Journal of Political Economy*, 126, 1219–1249.
- MARSH, T. A. AND R. C. MERTON (1987): "Dividend Behavior for the Aggregate Stock Market," *The Journal of Business*, 60, 1–40.
- MCLACHLAN, G. AND D. PEEL (2004): *Finite Mixture Models*, New York: John Wiley & Sons.
- MINCER, J. (1958): "Investment in Human Capital and Personal Income Distribution," *Journal of Political Economy*, 66, 281–302.
- (1974): *Schooling, Experience and Earnings*, New York: National Bureau of Economic Research.
- MINCER, J. AND B. JOVANOVIĆ (1981): "Labor Mobility and Wages," in *Studies in Labor Markets*, ed. by S. Rosen, Chicago: University of Chicago Press, chap. 5, 21–64.
- MÜLLER, H. G. (2012): *Nonparametric Regression Analysis of Longitudinal Data*, vol. 46, New York: Springer-Verlag.
- MÜLLER, H.-G., R. SEN, AND U. STADTMÜLLER (2011): "Functional Data Analysis for Volatility," *Journal of Econometrics*, 165, 233–245.
- MURPHY, K. M. AND F. WELCH (1990): "Empirical Age-Earnings Profiles," *Journal of Labor Economics*, 8, 202–229.
- NELSEN, R. (2007): *An Introduction to Copulas*, New York: Springer-Verlag.
- NEWBY, W. K. (1991): "Uniform Convergence in Probability and Stochastic Equicontinuity," *Econometrica*, 59, 1161–1167.
- NEWBY, W. K. AND D. MCFADDEN (1994): "Large Sample Estimation and Hypothesis Testing," in *Handbook of Econometrics*, ed. by R. F. Engle and D. McFadden, Amsterdam: Elsevier, vol. 4 of *Handbook of Econometrics*, chap. 36, 2111–2245.
- NEYMAN, J. (1959): "Optimal Asymptotic Tests of Composite Statistical Hypotheses," in *Probability and Statistics, The Harald Cramér Volume*, ed. by U. Grenander, New York: Wiley, 213–234.
- NING, W., A. K. GUPTA, C. YU, AND S. ZHANG (2009): "A Moment-Based Test for Homogeneity in Finite Mixture Models," *Communications in Statistics-Theory and Methods*, 38, 1371–1382.
- NIU, X., P. LI, AND P. ZHANG (2011): "Testing Homogeneity in a Multivariate Mixture Model," *The Canadian Journal of Statistics*, 39, 218–238.
- PEARSON, E. S. (1950): "On Questions Raised by the Combination of Tests Based on Discontinuous Distributions," *Biometrika*, 37, 383–398.
- PETERSEN, A. AND H.-G. MÜLLER (2016): "Functional Data Analysis for Density Functions by Transformation to a Hilbert Space," *The Annals of Statistics*, 44, 183 – 218.



- POLLARD, D. (1980): "The Minimum Distance Method of Testing," *Metrika*, 27, 43–70.
- RAMANATHAN, T. V. AND M. B. RAJARSHI (1992): "Rank Tests for Testing Randomness of a Regression Coefficient in a Linear Regression Model," *Metrika*, 39, 113–124.
- RAMSAY, J. AND B. SILVERMAN (1997): *The Analysis of Functional Data*, New York: Springer.
- (2002): *Applied Functional Data Analysis: Methods and Case Studies*, New York: Springer.
- RAMSAY, J. O. AND C. J. DALZELL (1991): "Some Tools for Functional Data Analysis," *Journal of the Royal Statistical Society. Series B (Methodological)*, 53, 539–572.
- RICE, J. A. AND B. W. SILVERMAN (1991): "Estimating the Mean and Covariance Structure Nonparametrically When the Data are Curves," *Journal of the Royal Statistical Society. Series B (Methodological)*, 53, 233–243.
- ROSENBERG, B. (1973): "The Analysis of a Cross-Section of Time Series by Stochastically Convergent Parameter Regression," *Annals of Economic and Social Measurement*, 2, 399–428.
- SCHLATTMANN, P. (2009): *Medical Applications of Finite Mixture Models*, Berlin Heidelberg: Springer-Verlag.
- SEO, W.-K. AND B. K. BEARE (2019): "Cointegrated Linear Processes in Bayes Hilbert Space," *Statistics & Probability Letters*, 147, 90–95.
- SIMONSOHN, U., L. D. NELSON, AND J. P. SIMMONS (2014): "P-Curve: A Key to the File-Drawer," *Journal of Experimental Psychology: General*, 143, 534–547.
- SKLAR, A. (1959): "Fonctions de Répartition à n Dimensions et Leurs Marges," *Publications de l'Institut de Statistique de l'Université de Paris*, 8, 229–231.
- STINCHCOMBE, M. B. AND H. WHITE (1992): "Some Measurability Results for Extrema of Random Functions Over Random Sets," *The Review of Economic Studies*, 59, 495–514.
- (1998): "Consistent Specification Testing with Nuisance Parameters Present Only under the Alternative," *Econometric Theory*, 14, 295–325.
- STOCK, J. AND M. WATSON (1998): "Median Unbiased Estimation of Coefficient Variance in a Time-Varying Parameter Model," *Journal of the American Statistical Association*, 93, 349–358.
- SU, L., Z. SHI, AND P. C. B. PHILLIPS (2016): "Identifying Latent Structures in Panel Data," *Econometrica*, 84, 2215–2264.
- SWAMY, P. A. V. B. AND G. S. TAVLAS (1995): "Random Coefficient Models: Theory and Applications," *Journal of Economic Surveys*, 9, 165–196.
- SWAMY, P. A. V. B. AND P. A. TINSLEY (1980): "Linear Prediction and Estimation Methods for Regression Models with Stationary Stochastic Coefficients," *Journal of Econometrics*, 12, 103–142.
- VAN ZWET, W. R. AND J. OOSTERHOFF (1967): "On the Combination of Independent Test Statistics," *The Annals of Mathematical Statistics*, 38, 659–680.

- VUONG, Q. H. (1989): "Likelihood Ratio Tests for Model Selection and Non-Nested Hypotheses," *Econometrica*, 57, 307–333.
- WALD, A. (1943): "Tests of Statistical Hypotheses Concerning Several Parameters When the Number of Observations is Large," *Transactions of the American Mathematical Society*, 54, 426–482.
- WANG, Y. (2011): *Smoothing Splines: Methods and Applications*, New York: Chapman & Hall/CRC.
- WESTBERG, M. (1985): "Combining Independent Statistical Tests," *Journal of the Royal Statistical Society. Series D (The Statistician)*, 34, 287–296.
- WHITE, H. (1982): "Maximum Likelihood Estimation of Misspecified Models," *Econometrica*, 50, 1–25.
- (1994): *Estimation, Inference and Specification Analysis*, New York: Cambridge University Press.
- WISWALL, M. AND B. ZAFAR (2017): "Preference for the Workplace, Investment in Human Capital, and Gender," *The Quarterly Journal of Economics*, 133, 457–507.
- WONG, T. S. T. AND W. K. LI (2014): "Test for Homogeneity in Gamma Mixture Models Using Likelihood Ratio," *Computational Statistics & Data Analysis*, 70, 127–137.
- WOOLDRIDGE, J. M. (2010): *Econometric Analysis of Cross Section and Panel Data*, vol. 1, MIT Press Books. The MIT Press.
- ZHANG, J.-T. AND J. CHEN (2007): "Statistical Inferences for Functional Data," *The Annals of Statistics*, 35, 1052 – 1079.
- ZIMMER, D. M. (2012): "The Role of Copulas in the Housing Crisis," *The Review of Economics and Statistics*, 94, 607–620.

		Random Effect Model Estimates			FLS Estimates			
		$\hat{\theta}_{1n}$	$\hat{\theta}_{2n}$	$\hat{\theta}_{3n}$	$\hat{\theta}_{1n}$	$\hat{\theta}_{2n}$	$\hat{\theta}_{3n}$	
Quadratic	$n = 100$	bias	-0.1559	0.0100	0.0001	-0.0020	0.0003	0.0001
		RMSE	0.1616	0.0113	0.0001	0.0354	0.0046	0.0001
		MAPE	1.7855	14.8004	34.3178	0.3234	5.3394	32.3014
	$n = 300$	bias	-0.0979	0.0037	0.0001	-0.0012	0.0001	0.0001
		RMSE	0.1010	0.0048	0.0001	0.0210	0.0027	0.0001
		MAPE	1.1210	5.8866	19.9489	0.1915	3.1593	19.0344
	$n = 500$	bias	-0.0862	0.0024	0.0001	-0.0011	0.0001	0.0001
		RMSE	0.0883	0.0034	0.0001	0.0163	0.0021	0.0001
		MAPE	0.9866	4.0821	15.4756	0.1491	2.4345	14.7696
Cubic	$n = 100$	bias	-0.2368	0.0212	-0.0002	-0.0038	-0.0009	0.0001
		RMSE	0.2426	0.0247	0.0008	0.0389	0.0104	0.0007
		MAPE	2.7866	14.8579	10.3461	0.3663	5.6874	8.9748
	$n = 300$	bias	-0.1807	0.0156	-0.0003	-0.0025	-0.0012	0.0001
		RMSE	0.1832	0.0172	0.0005	0.0228	0.0062	0.0004
		MAPE	2.1261	10.7005	6.7877	0.2139	3.3585	5.3161
	$n = 500$	bias	-0.1694	0.0144	-0.0003	-0.0033	-0.0011	0.0001
		RMSE	0.1711	0.0155	0.0005	0.0188	0.0049	0.0003
		MAPE	1.9935	9.8807	5.9050	0.1779	2.7076	4.2286
Quartic	$n = 100$	bias	-0.4919	0.0732	-0.0031	-0.4159	0.0566	-0.0021
		RMSE	0.4957	0.0774	0.0041	0.4195	0.0612	0.0032
		MAPE	6.1361	18.9974	10.4001	5.1872	14.6840	8.0462
	$n = 300$	bias	-0.4401	0.0691	-0.0033	-0.4146	0.0562	-0.0020
		RMSE	0.4416	0.0706	0.0036	0.4159	0.0058	0.0025
		MAPE	5.5147	18.0797	10.2448	5.1951	14.7054	6.6191
	$n = 500$	bias	-0.4299	0.0684	-0.0033	-0.4145	0.0562	-0.0020
		RMSE	0.4308	0.0693	0.0035	0.4152	0.0572	0.0023
		MAPE	5.3926	17.9323	10.4133	5.1982	14.7399	6.4429
Quartic (r)	$n = 100$	bias	-0.1166	0.0070	-0.0001	-0.0024	0.0003	0.0001
		RMSE	0.1231	0.0076	0.0001	0.0357	0.0030	0.0001
		MAPE	1.2665	20.7027	47.2998	0.3064	7.3860	39.6026
	$n = 300$	bias	-0.0604	0.0157	0.0001	-0.0018	0.0001	0.0001
		RMSE	0.0646	0.0024	0.0001	0.0213	0.0018	0.0001
		MAPE	0.6640	5.8626	24.0223	0.1812	4.3530	23.8094
	$n = 500$	bias	-0.0491	0.0047	0.0001	-0.0019	0.0001	0.0001
		RMSE	0.0521	0.0014	0.0001	0.0169	0.0014	0.0001
		MAPE	0.5373	3.5254	17.8046	0.1461	3.3860	18.3335

Table 5.1: ACCURACY ASSESSMENT FOR MODEL ESTIMATES UNDER QUADRATIC, CUBIC, QUARTIC AND RESTRICTED QUARTIC MODEL SPECIFICATIONS. This table shows the accuracy assessments using the finite sample bias, root mean square error (RMSE) and mean absolute percentage error (MAPE) of the random effect and FLS estimations. The estimation errors of  $\hat{\theta}_{4n}$  and  $\hat{\theta}_{5n}$  in the cubic and quartic specifications are very small in both panel estimates and FLS estimates, and they are omitted for brevity. DGP:  $w_i(t) = \theta_{1*} + \theta_{2*}t + \theta_{3*}t^2 + \theta_{4*}t^3 + \theta_{5*}t^4 + \varepsilon_i(t)$  and  $d\varepsilon_i(t) = -\kappa_{\dagger}\varepsilon_i(t)dt + \sigma_{\dagger}d\mathcal{W}_i(t)$ . Data for the random effect model estimation:  $\{w_{it} : i = 1, 2, \dots, n \text{ and } t = 1, 2, \dots, T\}$ . Data for FLS estimation:  $\{g_i(\cdot) : [0, 40] \mapsto \mathbb{R} : i = 1, 2, \dots, n\}$ , with  $g_i(\cdot)$  being constructed by the local polynomial kernel method. Number of Experiments: 5,000.

Size of the Test Statistics						
Statistics	Levels \ n	25	50	100	300	500
$\mathcal{W}_n^b$	1%	6.36	3.29	2.10	1.38	1.15
	5%	12.26	8.30	7.01	5.25	5.45
	10%	17.50	13.71	12.12	10.48	10.51
$\mathcal{LM}_n^b$	1%	6.36	3.29	2.10	1.38	1.15
	5%	12.26	8.30	7.01	5.25	5.45
	10%	17.50	13.71	12.12	10.48	10.51
$\mathcal{QLR}_n^b$	1%	3.84	2.27	1.82	1.01	1.27
	5%	9.38	7.13	6.66	5.21	5.48
	10%	14.60	12.55	12.13	9.76	10.91
$C(\alpha)$	1%	0.86	1.00	1.12	1.02	1.14
	5%	4.78	5.30	5.10	5.34	5.12
	10%	10.10	10.24	10.52	10.50	10.72

Power of the Test Statistics (Level of Significance: 5%)						
Statistics	$\pi_{\dagger} \setminus n$	25	50	100	300	500
$\mathcal{W}_n^b$	0.10	13.76	9.56	10.16	13.46	20.20
	0.20	16.20	17.02	19.86	43.28	65.94
	0.30	20.42	24.78	37.38	79.32	96.04
	0.40	29.08	38.96	61.68	97.38	99.86
	0.50	37.08	54.30	80.58	99.94	100.0
$\mathcal{LM}_n^b$	0.10	13.76	9.56	10.16	13.46	20.20
	0.20	16.20	17.02	19.86	43.28	65.94
	0.30	20.42	24.78	37.38	79.32	96.04
	0.40	29.08	38.96	61.68	97.38	99.86
	0.50	37.08	54.30	80.58	99.94	100.0
$\mathcal{QLR}_n^b$	0.10	9.34	8.84	10.40	16.84	25.42
	0.20	12.94	16.52	22.76	53.38	75.22
	0.30	17.88	26.76	44.34	86.88	98.38
	0.40	27.06	44.34	71.06	99.10	100.0
	0.50	37.34	61.82	87.70	99.96	100.0
$C(\alpha)$	0.10	8.10	10.10	13.46	26.22	36.32
	0.20	13.32	20.04	31.44	63.80	82.60
	0.30	21.40	31.78	54.12	90.74	98.72
	0.40	29.76	49.82	76.24	99.16	100.0
	0.50	40.10	64.92	90.20	99.94	100.0

Table 5.2: SIZE AND POWER OF THE TEST STATISTICS FOR HOMOGENEITY (IN PERCENT). This table shows the empirical rejection rates of the Wald, LM, QLR, and  $C(\alpha)$  test statistics. Null DGP:  $x_i \sim \text{IID Exp}(1)$ . Alternative DGP:  $x_i \sim \text{IID } \pi_{\dagger} \text{Exp}(1) + (1 - \pi_{\dagger}) \text{Exp}(2)$ . Model:  $\rho(\gamma, \theta_1, \theta_2) = \theta_1 + \theta_2(\gamma - 1)/(2\gamma - 1)^{1/2}$ ,  $\gamma \in [1.5, 2.5]$ . Number of Experiments: 5,000.

Size of the Test Statistics						
Statistics	Levels \ n	25	50	100	300	500
$\mathcal{W}_n^\sharp$	1%	3.42	3.26	2.86	2.75	1.52
	5%	10.20	9.92	8.42	6.15	5.93
	10%	19.88	17.85	15.53	14.35	11.10
$\mathcal{LM}_n^\sharp$	1%	1.18	1.98	2.38	2.65	1.61
	5%	8.58	8.52	8.27	6.22	5.71
	10%	17.24	17.54	15.24	14.20	10.79
$\mathcal{QLR}_n^\sharp$	1%	1.72	1.70	1.14	1.10	1.01
	5%	4.58	4.98	4.96	5.15	4.82
	10%	8.96	10.40	10.26	11.20	9.45
$\mathcal{BP}_n$	1%	0.68	0.80	0.88	1.02	1.12
	5%	3.02	4.22	4.50	4.90	4.94
	10%	6.78	8.52	9.66	10.04	9.12

Power of the Test Statistics (Level of Significance: 5%)						
Statistics	$\pi_* \setminus n$	25	50	100	300	500
$\mathcal{W}_n^\sharp$	0.01	12.00	10.22	8.40	6.18	6.40
	0.02	12.18	10.29	8.58	7.90	11.36
	0.03	12.46	10.62	9.14	12.26	20.80
	0.04	12.82	11.62	10.72	18.46	32.20
	0.05	13.14	12.52	12.48	26.42	43.78
$\mathcal{LM}_n^\sharp$	0.01	8.68	9.06	8.00	6.02	6.30
	0.02	8.72	9.15	8.44	7.68	11.12
	0.03	9.16	9.30	8.46	11.76	20.52
	0.04	9.74	10.08	9.82	18.10	31.82
	0.05	9.90	10.86	11.64	26.00	43.36
$\mathcal{QLR}_n^\sharp$	0.01	6.12	6.04	6.64	7.16	8.42
	0.02	7.12	7.62	9.56	12.80	17.92
	0.03	8.30	10.04	12.76	20.52	31.38
	0.04	9.42	12.12	16.00	29.88	44.86
	0.05	10.68	14.46	19.70	39.74	57.50
$\mathcal{BP}_n$	0.01	3.60	5.10	5.58	5.89	6.30
	0.02	4.16	5.76	6.70	7.10	7.12
	0.03	4.72	6.52	8.06	9.18	9.38
	0.04	5.28	7.48	8.96	10.44	10.64
	0.05	5.76	8.24	9.78	11.54	11.92

Table 5.3: SIZE AND POWER OF THE RANDOM COEFFICIENT TESTS (IN PERCENT). This table shows the empirical rejection rates of the Wald, LM, QLR, and Breusch-Pagan ( $\mathcal{BP}_n$ ) test statistics. Null DGP:  $y_i = x_i' \psi_\dagger + u_i$  such that  $x_i := (1, z_i)'$  and  $z_i \sim \text{IID } U[0, 1]$  and  $u_i \sim \text{IID } \mathcal{N}(0, 1)$ , where  $\psi_\dagger = (1, 1)'$ . Alternative DGP:  $y_i = x_i' \psi_\dagger + u_i$  such that  $u_i = \pi_\dagger^{1/2} x_i' \Omega^{1/2} (\gamma_\dagger) \nu_i + \delta_\dagger^{1/2} \varepsilon_i$  such that  $x_i := (1, z_i)'$ ,  $z_i \sim \text{IID } U[0, 1]$ , and  $(\nu_i', \varepsilon_i)' \sim \text{IID } \mathcal{N}(0, I_3)$ , where  $\psi_\dagger = (1, 1)'$ ,  $\gamma_\dagger = 0.5$ , and  $\delta_\dagger = 1$ . Model:  $\rho(\gamma, \theta_1, \theta_2) = \theta_1 + \theta_2 \exp(\gamma)$ ,  $\gamma \in [0, 1]$ . Number of Experiments: 5,000.

Size of the Test Statistics						
Statistics	Levels \ $n$	25	50	100	300	500
$\mathcal{W}_n^\sharp$	1%	2.45	1.64	1.40	0.97	0.99
	5%	8.28	6.02	5.96	4.87	4.77
	10%	14.03	11.30	10.92	9.73	9.77
$\mathcal{LM}_n^\sharp$	1%	0.66	0.84	1.00	0.82	0.88
	5%	5.14	4.62	5.18	4.48	4.45
	10%	10.62	9.60	9.83	9.20	9.23
$\mathcal{QLR}_n^\sharp$	1%	4.43	2.76	2.14	1.55	1.25
	5%	11.50	7.86	7.35	6.19	5.79
	10%	17.15	13.38	12.82	11.43	11.31
$\mathcal{CM}_n$	1%	0.18	0.72	1.00	0.73	1.47
	5%	1.14	1.98	2.82	3.20	3.24
	10%	2.18	3.36	4.96	6.40	5.90
$\mathcal{KS}_n$	1%	0.00	0.00	0.06	0.07	0.29
	5%	0.10	0.16	0.24	0.67	0.88
	10%	0.20	0.36	0.72	1.53	1.77

Power of the Test Statistics (Level of Significance: 5%)						
Statistics	$\pi_\dagger \setminus n$	25	50	100	300	500
$\mathcal{W}_n^\sharp$	0.10	7.46	6.72	6.06	8.10	10.62
	0.20	8.58	8.62	9.20	18.04	27.58
	0.30	9.32	11.56	15.58	35.86	52.46
	0.40	11.70	15.16	24.52	55.62	77.20
	0.50	13.74	20.66	33.06	74.44	93.24
$\mathcal{LM}_n^\sharp$	0.10	4.56	5.24	5.42	7.56	10.20
	0.20	5.60	6.38	7.94	17.00	25.98
	0.30	6.32	9.38	13.66	33.80	49.42
	0.40	8.30	12.86	21.88	52.92	74.28
	0.50	9.62	17.20	30.58	71.38	88.84
$\mathcal{QLR}_n^\sharp$	0.10	10.34	8.98	9.00	9.78	12.02
	0.20	11.56	10.44	10.94	20.36	30.10
	0.30	12.80	14.18	18.32	39.02	55.56
	0.40	15.54	18.36	27.16	59.06	79.78
	0.50	17.76	23.68	36.58	77.66	94.42
$\mathcal{CM}_n$	0.10	2.30	3.43	5.03	6.20	9.50
	0.20	2.20	4.17	9.53	21.00	30.32
	0.30	3.20	6.70	13.87	34.40	48.21
	0.40	3.60	9.73	21.87	54.40	74.12
	0.50	5.10	14.00	32.27	75.00	91.12
$\mathcal{KS}_n$	0.10	0.10	0.47	0.87	0.60	2.08
	0.20	0.00	0.37	0.90	4.40	7.00
	0.30	0.50	0.60	1.93	7.20	15.18
	0.40	0.20	1.20	3.47	19.20	37.35
	0.50	1.10	2.03	6.13	34.80	59.47

Table 5.4: SIZE AND POWER OF THE INDEPENDENCE TESTS (IN PERCENT). This table shows the empirical rejection rates of the Wald, LM, QLR, Cramér-von Mises ( $\mathcal{CM}_n$ ), and Kolmogorov-Smirnov ( $\mathcal{KS}_n$ ) test statistics. DGP: The copula  $c(u, v; \pi_\dagger, \gamma_{1\dagger}, \gamma_{2\dagger}) := (1 - \pi_\dagger)c_1(u, v; \gamma_{1\dagger}) + \pi_\dagger c_2(u, v; \gamma_{2\dagger})$ , with two margins  $x_i \sim \text{IID } \mathcal{N}(0, 1)$  and  $y_i \sim \text{IID } \mathcal{N}(0, 5)$ . Here,  $c_1(\cdot)$  and  $c_2(\cdot)$  are the independence and FGM copulas, respectively. Model:  $\rho(\gamma, \theta_1, \theta_2) = \theta_1 + \theta_2\gamma$ ,  $\gamma \in [0, 1]$ . Number of Experiments: 5,000.

Estimated FMSE of the mean of the log income paths								
	Male				Female			
	Quadratic	Cubic	Quartic	Quartic (r)	Quadratic	Cubic	Quartic	Quartic (r)
w/o Degree	6.09	5.94	5.93	5.93	5.77	5.66	5.65	5.65
Bachelor	5.58	5.52	5.47	5.55	5.32	5.27	5.23	5.29
Master	5.55	5.48	5.42	5.57	5.09	5.03	4.97	5.09
Ph.D	5.57	5.48	5.41	5.59	5.50	5.45	5.40	5.52

Estimated FMSE of the scaled mean of the log income paths								
	Male				Female			
	Quadratic	Cubic	Quartic	Quartic (r)	Quadratic	Cubic	Quartic	Quartic (r)
w/o Degree	3.20	2.92	2.90	2.90	3.03	2.81	2.78	2.78
Bachelor	3.15	3.04	2.96	3.10	3.10	3.01	2.94	3.04
Master	3.34	3.21	3.12	3.37	3.17	3.08	2.99	3.18
Ph.D	3.30	3.14	3.02	3.34	3.49	3.41	3.32	3.51

Table 6.5: FUNCTIONAL MEAN SQUARED ERRORS (FMSE) USING FUNCTION DATA OVER 0 TO 40 WORK EXPERIENCE YEARS. This table shows the estimated FMSEs of the log income paths under the quadratic, cubic, quartic and restricted quartic specifications for each group of the workers classified according to their education levels and genders. Here, the restricted quartic specification means the model specified by Murphy and Welch (1990) given in (11).

Inference results on the mean of log income paths across different genders				
	Education Level	Wald	LM	QLR
Quadratic	w/o Degree	5.08	11.45**	11.51
	Bachelor	140.41**	124.26**	1426.65**
	Master	54.40**	29.02**	345.38**
	Ph.D	40.32**	42.77**	459.46**
Cubic	w/o Degree	20.02**	41.57**	28.50
	Bachelor	174.84**	130.21**	1437.72**
	Master	40.69**	32.37**	364.48**
	Ph.D	60.51**	48.23**	468.17**
Quartic	w/o Degree	24.88**	42.31**	30.16
	Bachelor	289.97**	153.36**	1439.01**
	Master	63.05**	34.45**	365.54**
	Ph.D	83.61**	55.99**	470.26**
Quartic (r)	w/o Degree	13.32**	19.73**	18.17
	Bachelor	150.16**	132.58**	1414.82**
	Master	33.72**	30.13**	352.93**
	Ph.D	43.57**	47.34**	457.77**

Inference results on the mean of the scaled log income paths across different genders				
	Education Level	Wald	LM	QLR
Quadratic	w/o Degree	17.23**	16.55**	9.94
	Bachelor	21.46**	41.00**	73.40**
	Master	6.65	13.33**	19.12
	Ph.D	0.71	1.17	3.14
Cubic	w/o Degree	13.53**	27.18**	26.89*
	Bachelor	26.58**	54.48**	84.47**
	Master	7.69	18.07**	38.22*
	Ph.D	6.08	13.14*	11.86
Quartic	w/o Degree	26.15**	29.77**	28.55*
	Bachelor	53.80**	57.17**	85.76**
	Master	13.18*	17.58**	39.28*
	Ph.D	9.97	17.55**	13.95
Quartic (r)	w/o Degree	12.86**	16.57**	16.56
	Bachelor	24.93**	39.59**	61.57**
	Master	7.72	12.18**	26.68*
	Ph.D	0.63	0.83	1.45

Table 6.6: INFERENCE RESULTS USING FUNCTION DATA OVER 0 TO 40 WORK EXPERIENCE YEARS. This table shows the Wald, LM, and QLR test statistics and the inference results for the null hypothesis of equal mean of log income paths across different genders. The figures attached by ‘\*’ and ‘\*\*’ indicate the rejection of the null hypothesis at the 5% and 1% significance level, respectively.



Inference results on the mean of log income paths across different education levels

		Male			Female		
		Wald	LM	QLR	Wald	LM	QLR
Quadratic	w/o Degree vs. Bachelor	1908.46**	1374.39**	24803.82**	1588.52**	1129.27**	16617.83**
	Bachelor vs. Master	385.68**	172.95**	2621.40**	300.29**	191.17**	2153.12**
	Master vs. Ph.D	56.11**	51.10**	333.48**	73.14**	42.08**	194.69**
Cubic	w/o Degree vs. Bachelor	2178.92**	1510.92**	25441.14**	2013.79**	1348.72**	16973.37**
	Bachelor vs. Master	234.10**	181.54**	280.39**	250.68**	197.24**	2215.04**
	Master vs. Ph.D	64.47**	48.02**	335.42**	50.23**	41.50**	203.51**
Quartic	w/o Degree vs. Bachelor	3398.34**	1585.75**	25512.52**	3021.71**	1495.33**	17010.54**
	Bachelor vs. Master	399.12**	204.82**	2806.21**	380.31**	229.96**	2215.59**
	Master vs. Ph.D	100.55**	58.26**	336.26**	78.14**	45.57**	203.62**
Quartic (r)	w/o Degree vs. Bachelor	1683.18**	1466.64**	24955.59**	1423.92**	1276.41**	16754.52**
	Bachelor vs. Master	192.97**	180.38**	2732.06**	216.33**	204.09**	2174.74**
	Master vs. Ph.D	62.83**	54.63**	333.22**	50.57**	42.72**	202.33**

Inference results on the scaled mean of log income paths across different education levels

		Male			Female		
		Wald	LM	QLR	Wald	LM	QLR
Quadratic	w/o Degree vs. Bachelor	255.01**	486.43**	637.70**	101.86**	216.82**	182.05**
	Bachelor vs. Master	44.21**	59.41**	6.77	24.62**	40.87**	47.08**
	Master vs. Ph.D	8.91*	9.75*	20.65	27.97**	26.33**	67.23**
Cubic	w/o Degree vs. Bachelor	420.04**	966.72**	1274.81**	207.00**	538.63**	537.59**
	Bachelor vs. Master	36.63**	86.33**	189.23**	26.99**	55.92**	109.00*
	Master vs. Ph.D	9.29	16.25**	22.59	19.03**	26.89**	76.05**
Quartic	w/o Degree vs. Bachelor	719.97**	946.54**	1346.18**	332.41**	541.68**	574.77**
	Bachelor vs. Master	87.17**	81.95**	191.55**	58.52**	60.14**	109.55**
	Master vs. Ph.D	17.85**	18.36**	23.43	28.35**	27.27**	76.16**
Quartic (r)	w/o Degree vs. Bachelor	381.66**	601.31**	789.26**	159.10**	247.77**	318.74**
	Bachelor vs. Master	36.31**	53.38**	117.40**	24.85**	37.37**	68.70**
	Master vs. Ph.D	9.30*	9.67*	20.39*	18.34**	26.08**	74.87**

Table 6.7: INFERENCE RESULTS USING FUNCTION DATA OVER 0 TO 40 WORK EXPERIENCE YEARS. This table shows the Wald, LM, and QLR test statistics and the inference results for the null hypothesis of equal mean of log income paths across different education levels. The figures attached by ‘\*’ and ‘\*\*’ indicate the rejection of the null hypothesis at the 5% and 1% significance level, respectively.

Estimated FMSE of the mean of the log income paths

	Male				Female			
	Quadratic	Cubic	Quartic	Quartic (r)	Quadratic	Cubic	Quartic	Quartic (r)
w/o Degree	5.17	5.14	5.13	5.13	4.84	4.81	4.81	4.80
Bachelor	4.74	4.71	4.70	4.70	4.48	4.45	4.45	4.44
Master	4.72	4.69	4.69	4.68	4.31	4.28	4.28	4.27
Ph.D	4.75	4.71	4.71	4.71	4.63	4.61	4.61	4.60

Estimated FMSE of the scaled mean of the log income paths

	Male				Female			
	Quadratic	Cubic	Quartic	Quartic (r)	Quadratic	Cubic	Quartic	Quartic (r)
w/o Degree	2.32	2.25	2.24	2.24	2.28	2.22	2.21	2.20
Bachelor	2.27	2.20	2.20	2.19	2.30	2.23	2.23	2.22
Master	2.24	2.18	2.18	2.17	2.21	2.15	2.15	2.14
Ph.D	2.20	2.12	2.11	2.10	2.24	2.19	2.19	2.17

Table 6.8: FUNCTIONAL MEAN SQUARED ERRORS (FMSE) USING FUNCTION DATA OVER 10 TO 40 WORK EXPERIENCE YEARS. This table shows the estimated FMSEs of the log income paths under the quadratic, cubic, quartic and restricted quartic specifications for each group of the workers classified according to their education levels and genders. Here, the restricted quartic specification means the model specified by Murphy and Welch (1990) given in (11).

Inference results on the mean of log income paths across different genders				
	Education Level	Wald	LM	QLR
Quadratic	w/o Degree	1.20	1.49	6.85
	Bachelor	164.48 **	132.47**	1243.35**
	Master	37.62**	31.32**	316.42**
	Ph.D	34.80**	34.93**	347.86**
Cubic	w/o Degree	8.86	10.30	8.16
	Bachelor	435.48**	146.03**	1244.39**
	Master	98.55**	34.18**	316.65**
	Ph.D	131.26**	45.65**	349.11**
Quartic	w/o Degree	16.39**	12.12*	8.77
	Bachelor	592.76**	163.50**	1244.44**
	Master	131.90**	36.13**	316.73**
	Ph.D	199.89**	48.03**	349.29**
Quartic (r)	w/o Degree	1.51	1.44	7.26
	Bachelor	172.20**	147.59**	1243.31**
	Master	39.38**	34.26**	316.41**
	Ph.D	45.41**	45.79**	349.04**

Inference results on the mean of the scaled log income paths across different genders				
	Education Level	Wald	LM	QLR
Quadratic	w/o Degree	2.58	2.22	0.37
	Bachelor	0.57	1.35	2.22
	Master	0.12	0.19	0.41
	Ph.D	0.32	3.18	1.42
Cubic	w/o Degree	3.68	9.54*	1.68
	Bachelor	4.26	13.09*	3.27
	Master	0.84	2.65	0.65
	Ph.D	2.71	8.22	2.67
Quartic	w/o Degree	14.43*	11.03	2.28
	Bachelor	12.24*	14.65*	3.31
	Master	4.53	2.75	0.73
	Ph.D	10.65	7.98	2.85
Quartic (r)	w/o Degree	1.81	2.06	0.77
	Bachelor	0.59	1.29	2.18
	Master	0.10	0.15	0.41
	Ph.D	3.69	5.44	2.60

Table 6.9: INFERENCE RESULTS USING FUNCTION DATA OVER 10 TO 40 WORK EXPERIENCE YEARS. This table shows the Wald, LM, and QLR test statistics and the inference results for the null hypothesis of equal mean of log income paths across different genders. The figures attached by ‘\*’ and ‘\*\*’ indicate the rejection of the null hypothesis at the 5% and 1% significance level, respectively.

Inference results on the mean of log income paths across different education levels

		Male			Female		
		Wald	LM	QLR	Wald	LM	QLR
Quadratic	w/o Degree vs. Bachelor	1934.80**	1452.75**	21226.31**	1543.13**	1256.10**	13911.35**
	Bachelor vs. Master	218.02**	175.93**	2363.97**	235.96**	190.21**	1858.30**
	Master vs. Ph.D	52.66**	43.94**	287.84**	51.89**	39.87**	198.45**
Cubic	w/o Degree vs. Bachelor	6746.20**	1562.16**	21226.98**	5245.07**	1461.64**	13911.62**
	Bachelor vs. Master	822.53**	202.55**	2365.01**	760.05**	220.86**	1859.11**
	Master vs. Ph.D	147.68**	54.71**	288.56**	92.39**	42.28**	198.51**
Quartic	w/o Degree vs. Bachelor	8675.22**	1612.74**	21227.09**	6696.87**	1508.23**	13911.81**
	Bachelor vs. Master	966.42**	212.19**	2365.07**	960.79**	252.28**	1859.29**
	Master vs. Ph.D	157.25**	56.15**	288.69**	96.65**	43.91**	198.53**
Quartic (r)	w/o Degree vs. Bachelor	1916.51**	1536.87**	21226.46**	1629.18**	1412.39**	13911.53**
	Bachelor vs. Master	217.74**	192.23**	2364.33**	238.53**	210.17**	1858.61**
	Master vs. Ph.D	67.40**	52.16**	288.62**	53.45**	41.19**	198.55**

Inference results on the scaled mean of log income paths across different education levels

		Male			Female		
		Wald	LM	QLR	Wald	LM	QLR
Quadratic	w/o Degree vs. Bachelor	2.23	2.81	4.30	1.51	1.24	1.84
	Bachelor vs. Master	2.93	3.89	13.01	2.39	3.58	11.17
	Master vs. Ph.D	2.84	3.34	7.31	3.38	4.48	12.53
Cubic	w/o Degree vs. Bachelor	5.12	3.78	4.98	1.89	1.39	2.11
	Bachelor vs. Master	8.71	5.33	14.06	7.73	4.45	11.98
	Master vs. Ph.D	4.58	3.46	8.03	4.81	5.33	12.59
Quartic	w/o Degree vs. Bachelor	7.15	4.21	5.08	13.15*	7.12	2.30
	Bachelor vs. Master	19.39**	5.89	14.11	17.15**	4.52	12.16
	Master vs. Ph.D	11.11*	3.82	8.15	20.12**	6.41	12.61
Quartic (r)	w/o Degree vs. Bachelor	2.12	2.96	4.46	1.02	1.41	2.02
	Bachelor vs. Master	5.21	3.88	13.37	3.85	3.63	11.48
	Master vs. Ph.D	8.34*	3.94	8.09	2.96	5.06	12.63

Table 6.10: INFERENCE RESULTS USING FUNCTION DATA OVER 10 TO 40 WORK EXPERIENCE YEARS. This table shows the Wald, LM, and QLR test statistics and the inference results for the null hypothesis of equal mean of log income paths across different education levels. The figures attached by ‘\*’ and ‘\*\*’ indicate the rejection of the null hypothesis at the 5% and 1% significance level, respectively.

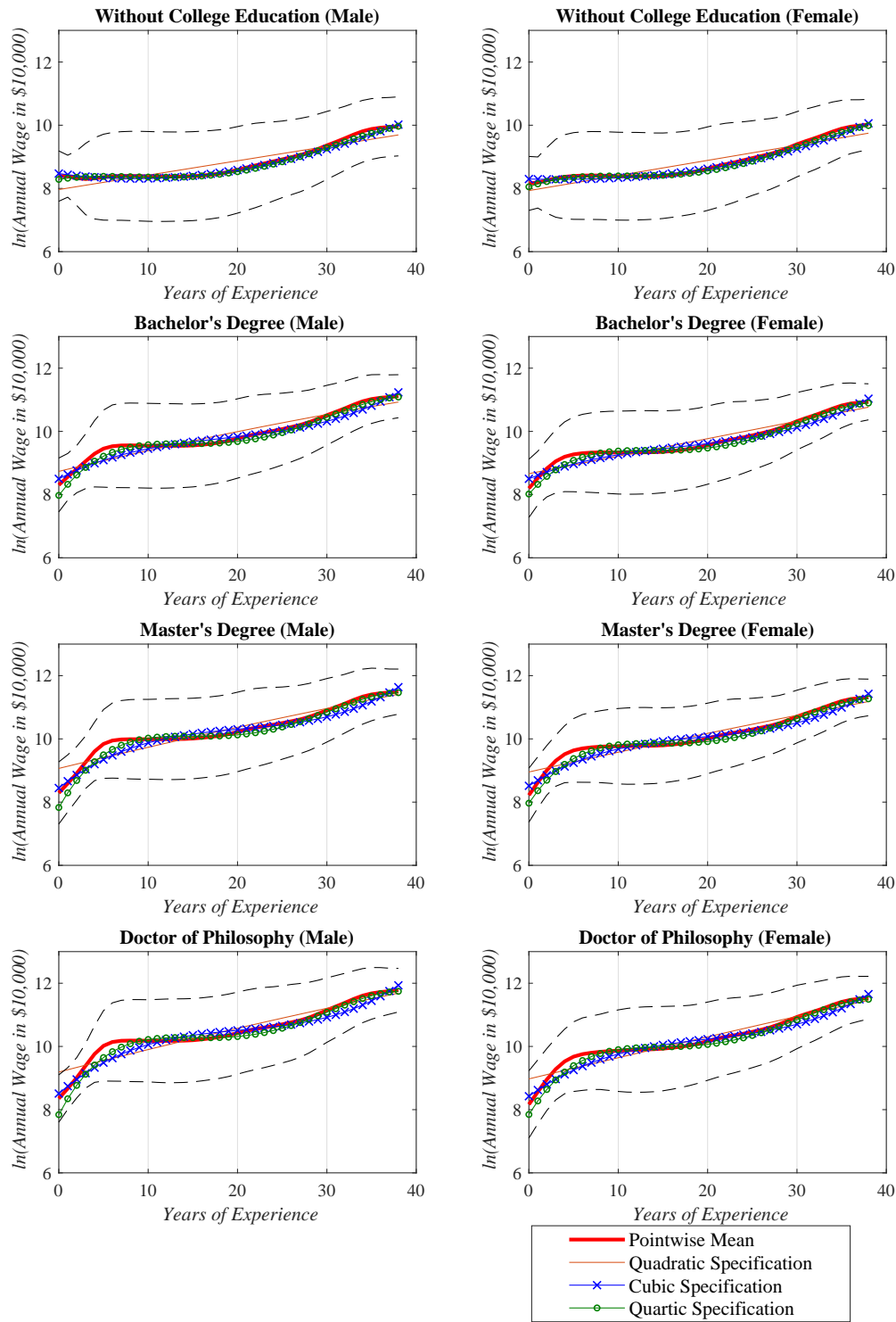


Figure 1: ESTIMATED MEAN CURVES OVER 0 TO 40 WORK EXPERIENCE YEARS. The red line corresponds to the pointwise mean of the individual log income paths, and the dotted lines correspond to its 80% bootstrap confidence bands. The mean estimates of log income paths under the quadratic, cubic and quartic specifications are displayed in brown, blue, and green lines.

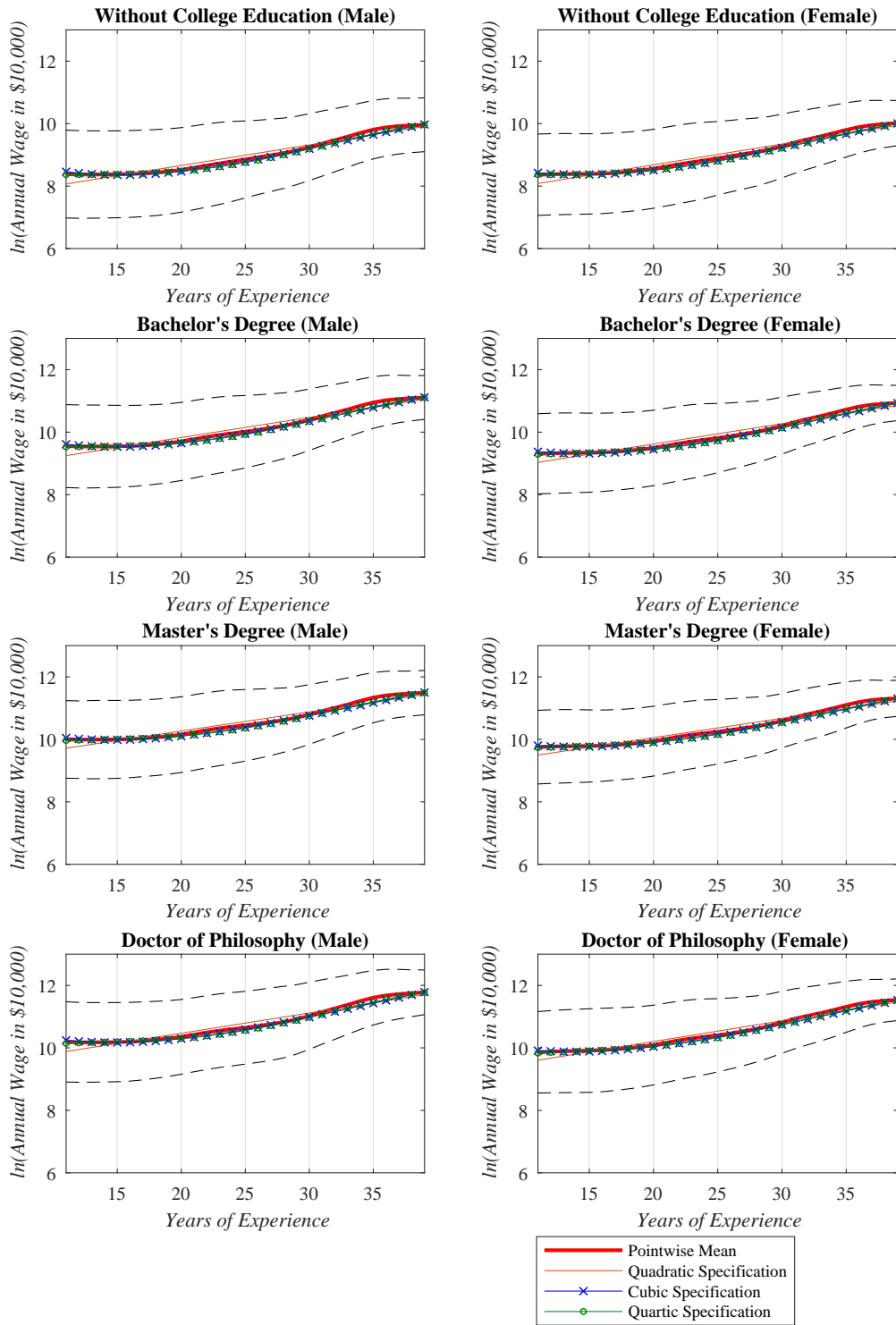


Figure 2: ESTIMATED MEAN CURVES OVER 10 TO 40 WORK EXPERIENCE YEARS. The red line corresponds to the pointwise mean of the individual log income paths, and the dotted lines correspond to its 80% bootstrap confidence bands. The mean estimates of log income paths under the quadratic, cubic and quartic specifications are displayed in brown, blue, and green lines.