

**ROBUSTNESS, INFINITESIMAL NEIGHBORHOODS,  
AND MOMENT RESTRICTIONS**

**By**

**Yuichi Kitamura, Taisuke Otsu and Kirill Evdokimov**

**August 2009**

**COWLES FOUNDATION DISCUSSION PAPER NO. 1720**



**COWLES FOUNDATION FOR RESEARCH IN ECONOMICS  
YALE UNIVERSITY  
Box 208281  
New Haven, Connecticut 06520-8281**

**<http://cowles.econ.yale.edu/>**

# ROBUSTNESS, INFINITESIMAL NEIGHBORHOODS, AND MOMENT RESTRICTIONS

YUICHI KITAMURA, TAISUKE OTSU, AND KIRILL EVDOKIMOV

ABSTRACT. This paper is concerned with robust estimation under moment restrictions. A moment restriction model is semiparametric and distribution-free, therefore it imposes mild assumptions. Yet it is reasonable to expect that the probability law of observations may have some deviations from the ideal distribution being modeled, due to various factors such as measurement errors. It is then sensible to seek an estimation procedure that are robust against slight perturbation in the probability measure that generates observations. This paper considers local deviations within shrinking topological neighborhoods to develop its large sample theory, so that both bias and variance matter asymptotically. The main result shows that there exists a computationally convenient estimator that achieves optimal minimax robust properties. It is semiparametrically efficient when the model assumption holds, and at the same time it enjoys desirable robust properties when it does not.

## 1. INTRODUCTION

Consider a probability measure  $P_0 \in \mathcal{M}$ , where  $\mathcal{M}$  is the set of all probability measures on the Borel  $\sigma$ -field  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$  of  $\mathcal{X} \subseteq \mathbb{R}^d$ . Let  $g : \mathcal{X} \times \Theta \rightarrow \mathbb{R}^m$  be a vector of functions parametrized by a  $p$ -dimensional vector  $\theta$  which resides in  $\Theta \subset \mathbb{R}^p$ . The function  $g$  satisfies:

$$(1.1) \quad E_{P_0} [g(x, \theta_0)] = \int g(x, \theta_0) dP_0 = 0, \quad \theta_0 \in \Theta.$$

The moment condition model (1.1) is semiparametric and distribution-free, therefore imposes mild assumptions. Nevertheless, it is reasonable to expect that the probability law of observations may have some deviations from the restriction under the moment condition model. It is then sensible to seek for

---

*Date:* This Version: June 4, 2009.

*Keywords:* Asymptotic Minimax Theorem, Hellinger Distance, Semiparametric Efficiency.

JEL Classification Number: C10.

We thank participants at the CIREQ Conference on GMM, the 2007 Winter Meetings of the Econometric Society, the 2007 Netherlands Econometrics Study Group Annual Conference, and seminars at Boston University, Chicago Booth, Harvard-MIT, LSE, Ohio State, Seoul National University, Texas A&M, the University of Tokyo, Vanderbilt, and Wisconsin for valuable comments. We acknowledge financial support from the National Science Foundation via grants SES-0241770 and SES-0551271 (Kitamura) and SES-0720961 (Otsu).

estimation and testing procedures that are robust against slight perturbations in the observed data, or more formally, perturbations in the probability measure that generates observations. This notion of robustness can be illustrated as follows. Let a functional  $\theta(P)$ ,  $P \in \mathcal{M}$  solve the moment condition model (1.1), in the sense that  $\theta_0 = \theta(P_0)$ . Suppose, however, observations  $x_1, \dots, x_n$  are not drawn according to  $P_0$ , but its “perturbed” version  $P$  instead. This can be attributed to various factors, including measurement errors or data contamination. These are imminent and realistic concerns in applications. The goal of robust estimation is to obtain an estimator  $\bar{\theta} = \bar{\theta}(x_1, \dots, x_n)$  that is not sensitive to such perturbations, so that the deviation of the estimated value  $\bar{\theta}$  from and the parameter value of interest  $\theta_0 = \theta(P_0)$  remains stable. Decompose the deviation as:

$$(1.2) \quad \bar{\theta} - \theta_0 = [\bar{\theta} - \theta(P)] + [\theta(P) - \theta(P_0)].$$

In the asymptotic MSE calculation presented below, the expectation of the square of the term in the first square bracket contributes to the variance of the estimator, whereas the second corresponds to the bias. An estimator that achieves small MSE *uniformly* in  $P$  over a neighborhood of  $P_0$  is desirable.

Asymptotic theory of robust estimation when the model is parametric has been considered extensively in the literature: see Rieder (1994) for a comprehensive survey. In a pioneering paper, Beran (1977) discusses “robust and efficient” estimation of parametric models. Suppose  $P_\theta, \theta \in \Theta \subset \mathbb{R}^k$  is a parametric family of probability measures. Observations are drawn from a probability law  $P$ , which may not be a member of the parametric family. Let  $p_\theta$  and  $p$  denote the densities associated with the probability measures  $P_\theta$  and  $P$ . It is well-known that the parametric MLE procedure corresponds to minimizing the objective function  $\rho = \int \log(p/p_\theta)p dx$ . Beran points out that a small change in the density  $p$  can lead to a large change in the objective function  $\rho$  (note the log in  $\rho$ ), implying the non-robustness of the MLE. He shows that the parametric Minimum Hellinger distance estimator (MHDE) is “robust and efficient,” in the sense that (i) it has an asymptotic minimax robust property and (ii) it is asymptotically efficient when the model assumption is satisfied, i.e. when the sample is generated from  $P_0 = P_{\theta_0}$ , where  $\theta_0$  is the true value of the parameter of interest. Let  $H(P_\theta, P) = \sqrt{\int (p_\theta^{1/2}(x) - p^{1/2}(x))^2 dx}$  denote the Hellinger distance between  $P_\theta$  and  $P$  (a slightly more general definition of the Hellinger distance is given in the next section). The MHDE for the parametric model is

$$\begin{aligned} \hat{\theta} &= \operatorname{argmin}_{\theta} H(P_\theta, \hat{P}) \\ &= \operatorname{argmin}_{\theta} \int (p_\theta^{1/2}(x) - \hat{p}^{1/2}(x))^2 dx \end{aligned}$$

where  $\hat{p}$  is a nonparametric density estimator, such as a kernel density estimator, for  $P$  and  $\hat{P}$  is the corresponding estimator for the probability measure of  $x$ . The MHDE is asymptotically equivalent to MLE and thus efficient if the model assumption is satisfied. One can replace the Hellinger distance with other divergence measures such as the Kolmogorov-Smirnov distance, which would make the corresponding minimum divergence estimator even more robust, but it would incur efficiency loss. The parametric MHDE has been studied extensively and applied to various models.

The parametric MHDE has theoretical advantages and excellent finite sample performance documented by numerous simulation studies, but it has limitations as well. It requires the nonparametric density estimator when at least some components of  $x$  are continuously distributed. This makes its practical application inconvenient, and is problematic when  $x$  is high dimensional, due to the curse of dimensionality. It also necessitates the evaluation of the integral  $\int (p_\theta^{1/2}(x) - \hat{p}^{1/2}(x))^2 dx$ . This would either involve numerical integration or an approximation by an empirical average with inverse density weighting using a nonparametric density estimator. The former can be hard to implement, and the latter may have undesirable effects in finite samples. This paper aims at developing robust methods for moment restriction models, by applying the MHDE procedure. The resulting estimator is semiparametrically efficient when the model assumption holds, and at the same time it enjoys an optimal minimax robust property when it does not. The implementation of the estimator is easy. Unlike its parametric predecessor, it requires neither nonparametric density estimation nor evaluation of integration.

## 2. PRELIMINARIES

Suppose a random sample  $\{x_i\}_{i=1}^n$  generated from  $P$  is observed. The econometrician wishes to estimate the unknown  $\theta_0$  in (1.1) from the sample. As discussed in Section 1, our focus is on robust estimation of  $\theta_0$  when the probability measure  $P$ , from which the observations are drawn, is a (locally) perturbed version of  $P_0$ , not  $P_0$  itself. There exists an extensive literature concerning the estimation of (1.1) under the “classical” setting where data are indeed drawn from  $P_0$ . Many estimators for  $\theta_0$  are available, including GMM (Hansen (1982)), the empirical likelihood (EL) estimator and its variants. This paper is concerned with an estimator that can be viewed as the MHDE applied to the moment restriction model (1.1). The Hellinger distance between two probability measures is defined as follows:

**Definition 2.1.** Let  $P$  and  $Q$  be probability measures with densities  $p$  and  $q$  with respect to a dominating measure  $\nu$ . The Hellinger distance between  $P$  and  $Q$  is then given by:

$$H(P, Q) = \left\{ \int (p^{1/2} - q^{1/2})^2 d\nu \right\}^{1/2} = \left\{ 2 - 2 \int p^{1/2} q^{1/2} d\nu \right\}^{1/2}.$$

The following notation that does not explicitly refer to the dominating measure is often convenient:

$$H(P, Q) = \left\{ \int (dP^{1/2} - dQ^{1/2})^2 \right\}^{1/2} = \left\{ 2 - 2 \int dP^{1/2} dQ^{1/2} \right\}^{1/2}.$$

We now provide some results concerning the Hellinger distance that are useful in understanding the robustness theorems in the next section.

**Definition 2.2.** Let  $P$  and  $Q$  be probability measures with densities  $p$  and  $q$  with respect to a dominating measure  $\nu$ . The  $\alpha$ -divergence from  $Q$  to  $P$  is given by

$$I_\alpha(P, Q) = \frac{1}{\alpha(1-\alpha)} \int \left( 1 - \left( \frac{p}{q} \right)^\alpha \right) q d\nu, \quad \alpha \in \mathbb{R}.$$

If  $P$  is not absolutely continuous respect to  $Q$ , then  $\int \mathbb{I}\{p > 0, q = 0\} d\nu > 0$ , and as a consequence  $I_\alpha(P, Q) = \infty$  for  $\alpha \geq 1$ . A similar argument shows that  $I_\alpha(P, Q) = \infty$  if  $Q \not\ll P$  and  $\alpha \leq 0$ . Note that  $I_\alpha$  is well-defined for  $\alpha = 1$  by taking the limit  $\alpha \rightarrow 1$  in the definition. Indeed, L'Hospital's Rule implies that

$$\lim_{\alpha \rightarrow 1} I_\alpha(P, Q) = \int \log \left( \frac{p}{q} \right) p d\nu := K(P, Q)$$

(with the above convention for the case where  $P \not\ll Q$ ), giving rise to the well-known Kullback-Leibler (KL) divergence measure from  $P$  to  $Q$ . The case with  $\alpha = 0$  corresponds to the KL divergence with the roles of  $P$  and  $Q$  reversed. Note that the  $\alpha$ -divergence includes the Hellinger distance as a special case, in the sense that

$$H^2(P, Q) = \frac{1}{2} I_{\frac{1}{2}}(P, Q).$$

The following Lemma provides an upper bound for the Hellinger distance. It generalizes well-known information theoretic inequalities.

**Lemma 2.1.** For probability measures  $P$  and  $Q$ ,

$$(2.1) \quad \max(\alpha, 1 - \alpha) I_\alpha(P, Q) \geq \frac{1}{2} I_{\frac{1}{2}}(P, Q)$$

for every  $\alpha \in \mathbb{R}$ .

*Proof.* We first show the claim for  $\alpha < \frac{1}{2}$ , that is,

$$(2.2) \quad (1 - \alpha) I_\alpha(P, Q) - \frac{1}{2} I_{\frac{1}{2}}(P, Q) \geq 0.$$

Let  $H_\alpha(x) = \frac{1}{\alpha}(1 - x^\alpha) - 2\left(1 - x^{\frac{1}{2}}\right)$ ,  $0 \leq x \leq \infty$ , then the above inequality becomes

$$(2.3) \quad \int H_\alpha\left(\frac{p}{q}\right) q d\nu \geq 0.$$

Note

$$\frac{d}{dx} H_\alpha(x) = -x^{\alpha-1} + x^{-\frac{1}{2}} \begin{cases} > 0 & \text{if } x > 1 \\ = 0 & \text{if } x = 1 \\ < 0 & \text{if } x < 1. \end{cases}$$

The above holds for the case with  $\alpha = 0$  as well, since  $H_0(x) = -\log x - 2\left(1 - x^{\frac{1}{2}}\right)$ . Moreover,  $H_\alpha(1) = 0$ . Therefore  $H_\alpha(x) \geq 0$  for all  $x \geq 0$ , and the desired inequality (2.3) follows immediately. Next, we prove the case with  $\alpha > \frac{1}{2}$ , that is,

$$\alpha I_\alpha(P, Q) \geq \frac{1}{2} I_{\frac{1}{2}}(P, Q).$$

Let  $\beta = 1 - \alpha < \frac{1}{2}$ , then the above inequality becomes

$$(2.4) \quad (1 - \beta) I_{1-\beta}(P, Q) \geq \frac{1}{2} I_{\frac{1}{2}}(P, Q).$$

By (2.2) and the symmetry of the Hellinger distance,

$$(1 - \beta) I_\beta(Q, P) \geq \frac{1}{2} I_{\frac{1}{2}}(Q, P) = \frac{1}{2} I_{\frac{1}{2}}(P, Q).$$

But the equality  $I_{1-\beta}(P, Q) = I_\beta(Q, P)$  holds for every  $\beta \in \mathbb{R}$ , and (2.4) follows.  $\square$

**Remark 2.1.** Lemma 2.1 has some implications on a neighborhood system generated by the Hellinger distance. Consider the following neighborhood of a probability measure  $P$  with its radius in terms of  $I_\alpha$  is  $\delta > 0$ :

$$B_{I_\alpha}(P, \delta) = \left\{ Q \in \mathcal{M} : \sqrt{I_\alpha(Q, P)} \leq \delta \right\}.$$

Lemma 2.1 implies that

$$I_\alpha(P, Q) \geq \frac{1}{2\left(\left(\frac{1}{2} + L\right) \vee \left(\frac{1}{2} + U\right)\right)} I_{\alpha_0}(P, Q)$$

holds for every  $\alpha \in \left[\frac{1}{2} - L, \frac{1}{2} + U\right]$  where  $L, U > 0$  determine the lower and upper bounds for the range of  $\alpha$ , if  $\alpha_0 = \frac{1}{2}$ . It is easy to verify that this statement holds only if  $\alpha_0 = \frac{1}{2}$ . Now, define

$$C(L, U) = \left(\frac{1}{2} + L\right) \vee \left(\frac{1}{2} + U\right),$$

then by the above inequality

$$(2.5) \quad \cup_{\alpha \in [\frac{1}{2} - L, \frac{1}{2} + U]} B_{I_\alpha}(P_0, \delta) \subset B_{I_{1/2}}\left(P_0, \sqrt{2C(L, U)}\delta\right).$$

That is, the union of the  $I_\alpha$ -based neighborhoods over  $\alpha \in [\frac{1}{2} - L, \frac{1}{2} + U]$  is covered by the Hellinger neighborhood  $B_{I_{1/2}}$  with a “margin” given by the multiplicative constant  $2\sqrt{C(L, U)}$ . (2.5) is important, since in what follows we consider robustness of estimators against perturbation of  $P_0$  within its neighborhood, and it is desirable to use a neighborhood that is sufficiently large to accommodate a large class of perturbations. The inclusion relationship shows that the Hellinger-based neighborhood covers other neighborhood systems based on  $I_\alpha, \alpha \in [\frac{1}{2} - L, \frac{1}{2} + U]$  if the radii are chosen appropriately. It is easy to verify that (2.5) does not hold if the Hellinger distance  $I_{\frac{1}{2}}$  is replaced by  $I_\alpha, \alpha \neq \frac{1}{2}$ , showing the special status of the Hellinger distance among the  $\alpha$ -divergence family.

**Remark 2.2.** Lemma 2.1 is a statement for every pair of measures  $(P, Q)$ , thus it holds even if  $P \not\ll Q$  or  $Q \not\ll P$ . On the other hand, it is useful to consider the behavior of  $I_\alpha$  when one of the two measures is not absolutely continuous with respect to the other. Consider a sequence of probability measures  $\{P^{(n)}\}_{i=1}^\infty$ . Suppose  $I_\alpha(P^{(n)}, P_0) \rightarrow 0$  for an  $\alpha \in \mathbb{R}$ , then  $I_{\alpha'}(P^{(n)}, P_0) \rightarrow 0$  for every  $\alpha' \in (0, 1)$ . But the reverse (i.e. reversing the roles of  $\alpha$  and  $\alpha'$ ) is not true. If  $P^{(n)}, n \in \mathbb{N}$  are not absolutely continuous respect to  $P_0$ ,  $I_{\alpha'}(P^{(n)}, P_0) = \infty$  for every  $\alpha' \geq 1$  even if  $\rho_\alpha(P^{(n)}, P_0) \rightarrow 0$  for  $\alpha \in (0, 1)$  (and a similar argument holds for  $\alpha' \leq 0$ ). This shows that  $I_\alpha$ -based neighborhoods with  $\alpha \notin (0, 1)$  are too small: there are measures that are outside of  $B_{I_\alpha}(P_0, \delta), \alpha \notin (0, 1)$  no matter how large  $\delta$  is, or how close they are to  $P_0$  in terms of, say, the Hellinger distance  $H$ . This shortcoming applies to neighborhoods based on the KL divergence and the  $\chi^2$  measure (see Remark 2.3), as they correspond to  $I_\alpha$  with  $\alpha = -1, 0, 1$  and 2.

**Remark 2.3.** The inequality in Lemma 2.1 might be of interest on its own as it generalizes various information theoretic inequalities in the literature. For  $\alpha = 1$  or 0, it corresponds to the well-known inequality between the KL divergence and the Hellinger distance

$$(2.6) \quad H(P, Q)^2 \leq K(P, Q),$$

see, for example, Pollard (2002), p.62. Also, consider the  $\chi^2$  distance, which is given by  $\chi^2(P, Q) = \int \frac{(p-q)^2}{q} d\nu$ . Then

$$(2.7) \quad H(P, Q)^2 \leq \chi^2(P, Q)$$

holds (Reiss (1989)). This is a special case of Lemma 2.1 with  $\alpha = -1$  and 2. Proposition 3.1 in Zhang (2006) shows that the inequality (2.1) holds for  $\alpha \in (0, 1)$ . Note that Zhang's result interpolates (2.6) and the same inequality with  $P$  and  $Q$  reversed, but covers neither (2.6) nor (2.7)<sup>1</sup>. These results have been obtained more or less on a case by case basis. Lemma 2.1 proves that this type of inequality holds for all  $\alpha \in \mathbb{R}$ , unifying those well-established results in the literature.

Beran (1977), considering a parametric model, proposed MHDE that minimizes the Hellinger distance between a model-based probability measure (from the parametric family) and a nonparametric probability measure estimate. An application of the MHDE procedure to the moment condition model (1.1) yields a computationally simple procedure as follows. Let  $P_n$  denote the empirical measure of observations  $\{x_i\}_{i=1}^n$ .  $P_n$  is an appropriate model-free estimator in our construction of the MHDE. Let

$$\mathcal{P}_\theta = \left\{ P \in \mathcal{M} : \int g(x, \theta) dP = 0 \right\}$$

and

$$(2.8) \quad \mathcal{P} = \cup_{\theta \in \Theta} \mathcal{P}_\theta,$$

then the MHDE, denoted by  $\hat{\theta}$ , is defined to be a parameter value that solves the optimization problem

$$\inf_{\theta \in \Theta} \inf_{P \in \mathcal{P}_\theta} H(P, P_n) = \inf_{P \in \mathcal{P}} H(P, P_n).$$

By convex duality theory (Kitamura (2006)), the objective function has the following representation:

$$\inf_{P \in \mathcal{P}_\theta} H(P, P_n) = \max_{\gamma \in \mathbb{R}^m} -\frac{1}{n} \sum_{i=1}^n \frac{1}{1 + \gamma' g(x_i, \theta)}$$

Therefore the MHDE is  $\hat{\theta} = \arg \min_{\theta \in \Theta} \max_{\gamma \in \mathbb{R}^m} -\frac{1}{n} \sum_{i=1}^n \frac{1}{1 + \gamma' g(x_i, \theta)}$ , which is easy to compute.

It is straightforward to verify that we can obtain the MHDE as a Generalized Empirical Likelihood (GEL) estimator by letting  $\gamma = -1/2$  in equation (2.6) of Newey and Smith (2004). Asymptotic properties of the (G)EL estimators for  $\theta_0$  in (1.1), when data drawn from  $P_0$  are observed are well-understood (see, for example, Kitamura and Stutzer (1997), Smith (1997), Imbens, Spady, and Johnson (1998), and Newey and Smith (2004)). Let  $G = E_{P_0} [\partial g(x, \theta_0) / \partial \theta']$ ,  $\Omega = E_{P_0} [g(x, \theta_0) g(x, \theta_0)']$ , and  $\Sigma = G' \Omega^{-1} G$ . Then

$$(2.9) \quad \sqrt{n} \left( \hat{\theta}_\alpha - \theta_0 \right) \xrightarrow{d} N(0, \Sigma^{-1}).$$

---

<sup>1</sup>Zhang (2006) also derives a lower bound for the Hellinger distance in terms of  $I_\alpha$ .



It follows that the MHDE and other GEL estimators are semiparametrically efficient in the absence of data perturbation. At the same time, the MHDE possesses a distinct property of being asymptotic optimal robust if observations are drawn from a perturbed version of  $P_0$ , as we shall see in the next section.

### 3. ROBUST ESTIMATION THEORY

We now analyze robustness of the MHDE  $\hat{\theta}$ . Define a functional

$$T(P) = \arg \min_{\theta \in \Theta} \max_{\gamma \in \mathbb{R}^m} - \int \frac{1}{1 + \gamma' g(x, \theta)} dP$$

then the MHDE can be interpreted as the value of functional  $T$  evaluated at the empirical measure  $P_n$ . In other words, each realization of  $P_n$  completely determines the value of the MHDE  $\hat{\theta}$ . To make the dependence explicit, we write  $\hat{\theta} = T(P_n)$ , and study properties of the mapping  $T : \mathcal{M} \rightarrow \Theta$ . This definition of  $T(\cdot)$ , however, causes a technical difficulty when the distribution of  $g(x, \theta)$  is unbounded for some  $\theta \in \Theta$  and  $P \in \mathcal{M}$ . To overcome this, we introduce the following mapping defined by a trimmed moment function:

$$\bar{T}(Q) = \arg \min_{\theta \in \Theta} \inf_{P \in \bar{\mathcal{P}}_\theta} H(P, Q),$$

where  $\{m_n\}_{n \in \mathbb{N}}$  is a sequence of positive numbers satisfying  $m_n \rightarrow \infty$  as  $n \rightarrow \infty$ , and

$$\begin{aligned} \bar{\mathcal{P}}_\theta &= \left\{ P \in \mathcal{M} : \int g(x, \theta) \mathbb{I}\{x \in \mathcal{X}_n\} dP = 0 \right\}, \\ \mathcal{X}_n &= \left\{ x \in \mathcal{X} : \sup_{\theta \in \Theta} |g(x, \theta)| \leq m_n \right\}, \end{aligned}$$

with the indicator function  $\mathbb{I}\{\cdot\}$  and the Euclidean norm  $|\cdot|$ , i.e.,  $\mathcal{X}_n$  is a trimming set to bound the moment function and  $\bar{\mathcal{P}}_\theta$  is a set of probability measures satisfying the bounded moment condition  $E_P[g(x, \theta) \mathbb{I}\{x \in \mathcal{X}_n\}] = 0$ . Lemma 7.1 (i) guarantees that for each  $n \in \mathbb{N}$  and  $Q \in \mathcal{M}$  the value  $\bar{T}(Q)$  exists.

Let  $\tau : \Theta \rightarrow \mathbb{R}$  be a possibly nonlinear transformation of the parameter. We consider the estimation problem of the transformed scalar parameter  $\tau(\theta_0)$ . The transformation  $\tau$  to a scalar, as used by Rieder (1994), is convenient in calculating squared biases and MSE's. One may, for example, choose  $\tau(\theta) = c'\theta$  using a  $p$ -vector  $c$ . We first investigate the behavior of the bias term  $\tau \circ \bar{T}(Q) - \tau(\theta_0)$  in a  $(\sqrt{n}$ -shrinking) Hellinger ball with radius  $r > 0$  around  $P_0$

$$B_H(P_0, r/\sqrt{n}) = \{Q \in \mathcal{M} : H(Q, P_0) \leq r/\sqrt{n}\}.$$

**Assumption 3.1.** *Suppose the following conditions hold:*

- (i):  $\{x_i\}_{i=1}^n$  is iid;
- (ii):  $\Theta$  is compact;
- (iii):  $\theta_0 \in \text{int}(\Theta)$  is a unique solution to  $E_{P_0} [g(x, \theta)] = 0$ ;
- (iv):  $g(x, \theta)$  is continuous over  $\Theta$  at each  $x \in \mathcal{X}$ ;
- (v):  $E_{P_0} [\sup_{\theta \in \Theta} |g(x, \theta)|^\eta] < \infty$  for some  $\eta > 2$ , and there exists a neighborhood  $\mathcal{N}$  around  $\theta_0$  such that  $E_{P_0} [\sup_{\theta \in \mathcal{N}} |g(x, \theta)|^4] < \infty$ ,  $g(x, \theta)$  is continuously differentiable a.s. on  $\mathcal{N}$ ,  $\sup_{x \in \mathcal{X}_n, \theta \in \mathcal{N}} |\partial g(x, \theta) / \partial \theta'| = o(n^{1/2})$ , and  $E_{P_0} [\sup_{\theta \in \mathcal{N}} |\partial g(x, \theta) / \partial \theta'|^2] < \infty$ ;
- (vi):  $G$  has the full column rank and  $\Omega$  is positive definite;
- (vii):  $\{m_n\}_{n \in \mathbb{N}}$  satisfies  $m_n \rightarrow \infty$ ,  $nm_n^{-\eta} \rightarrow 0$ , and  $n^{-1/2}m_n^{1+\epsilon} = O(1)$  for some  $0 < \epsilon < 2$  as  $n \rightarrow \infty$ ;
- (viii):  $\tau$  is continuously differentiable at  $\theta_0$ .

Assumptions 3.1 (i)-(vi) are standard in the literature on GMM. Assumption 3.1 (iii) is a global identification condition of the true parameter  $\theta_0$  under  $P_0$ . Assumption 3.1 (iv) ensures the continuity of the mapping  $\bar{T}(Q)$  in  $Q \in \mathcal{M}$  for each  $n \in \mathbb{N}$ . Assumption 3.1 (v) contains the smoothness and boundedness conditions for the moment function and its derivatives. This is stronger than a typical assumption imposed to obtain the standard asymptotic distribution (2.9). Assumption 3.1 (vi) is a local identification condition for  $\theta_0$ . This assumption guarantees that the asymptotic variance matrix  $\Sigma^{-1}$  exists. Assumption 3.1 (vii) is on the trimming parameter  $m_n$ . If  $m_n \sim n^a$ , this assumption is satisfied for  $1/\eta < a < 1/2$ . Assumption 3.1 (viii) is a standard requirement for the parameter transformation  $\tau$ .

To characterize a class of estimators to be compared with the MHDE, we introduce the following definition.

**Definition 3.1.** Let  $T_a(P_n)$  be an estimator of  $\theta_0$  based on a mapping  $T_a : \mathcal{M} \rightarrow \Theta$ . Also, let  $P_{\theta, \zeta}$  be a regular parametric submodel (see, Bickel, Klassen, Ritov, and Wellner (1993, p. 12)) of  $\mathcal{P}$  in (2.8) such that  $P_{\theta_0, 0} = P_0$ .

(i):  $T_a$  is called **Fisher consistent** if

$$(3.1) \quad \sqrt{n} \left( T_a \left( P_{\theta_0 + t/\sqrt{n}, \zeta_n} \right) - \theta_0 \right) \rightarrow t.$$

holds for every submodel  $P_{\theta, \zeta}$  that satisfies  $P_{\theta_0 + t/\sqrt{n}, \zeta_n} \in B_H(P_0, r/\sqrt{n})$  eventually with  $\zeta_n = O(n^{-1/2})$  and for every  $t \in \mathbb{R}^p$ .

(ii):  $T_a$  is called **regular** for  $\theta_0$  if for every  $\{P_{\theta_n, \zeta_n}\}_{n \in \mathbb{N}}$  with  $(\theta'_n, \zeta'_n)' = (\theta'_0, 0)' + O(n^{-1/2})$ , there exists a probability measure  $M$  such that

$$(3.2) \quad \sqrt{n}(T_a(P_n) - T_a(P_{\theta_n, \zeta_n})) \xrightarrow{d} M, \quad \text{under } P_{\theta_n, \zeta_n},$$

where the measure  $M$  does not depend on the sequence  $\{(\theta'_n, \zeta'_n)'\}_{n \in \mathbb{N}}$ .

Both conditions are weak and satisfied by GMM, (G)EL and other standard estimators. For example, the mapping  $T_a$  for the continuous updating GMM estimator (CUE) is given by

$$T_{CUE}(P) = \operatorname{argmin}_{\theta \in \Theta} \left[ \int g(x, \theta) dP \right]' \left[ \int g(x, \theta) g(x, \theta) dP \right]^{-1} \left[ \int g(x, \theta) dP \right],$$

and under Assumption 3.1  $T_{CUE}(P_{\theta_0 + t/\sqrt{n}, \zeta_n}) = \theta_0 + t/\sqrt{n}$  for large  $n$ . CUE therefore trivially satisfies (3.1). The regularity condition (3.2) is standard in the literature of semiparametric efficiency; see, for example, Bickel, Klassen, Ritov, and Wellner (1993) and Newey (1990).

The following theorem shows the optimal robustness of the (trimmed) MHDE in terms of its maximum bias.

**Theorem 3.1.** Suppose that Assumption 3.1 holds.

(i): For every  $T_a$  which is Fisher consistent,

$$\liminf_{n \rightarrow \infty} \sup_{Q \in B_H(P_0, r/\sqrt{n})} n(\tau \circ T_a(Q) - \tau(\theta_0))^2 \geq 4r^2 B^*,$$

for each  $r > 0$ , where  $B^* = \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' \Sigma^{-1} \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)$ .

(ii): The mapping  $\bar{T}$  is Fisher consistent and satisfies

$$\lim_{n \rightarrow \infty} \sup_{Q \in B_H(P_0, r/\sqrt{n})} n(\tau \circ \bar{T}(Q) - \tau(\theta_0))^2 = 4r^2 B^*,$$

for each  $r > 0$ .

**Remark 3.1.** The above result is concerned with deterministic properties of  $T_a$  and  $T$ .  $T_a(Q)$  and  $T(Q)$  can be regarded as the (probability) limit of the estimators  $T_a(P_n)$  and  $T(P_n)$  under  $Q$ , and therefore the terms evaluated here correspond to the bias of each estimators due to the deviation of  $Q$  from  $P_0$ . The theorem says that in the class of all mappings that are Fisher consistent, the mapping  $\bar{T}$  has the smallest maximum bias over the set  $B_H(P_0, r/\sqrt{n})$ . The (trimmed version of) the Hellinger-based mapping  $\bar{T}$  is therefore optimally robust asymptotically in a minimax sense. The term  $4r^2 B^*$  provides a sharp lower bound for maximum squared bias, and it is attained by  $\bar{T}$ .

**Remark 3.2.** The second part of the theorem deals with the trimmed version of the MHDE. It avoids the complications associated with the existence of  $T(Q)$  for certain  $Q$ 's. If the support of  $\sup_{\theta \in \Theta} |g(x, \theta)|$  is bounded under every  $Q \in B_H(P_0, r/\sqrt{n})$  for large enough  $n$  (e.g. if the moment function  $g$  is bounded), then we do not need the trimming term  $\mathbb{I}\{x \in \mathcal{X}_n\}$ . In this case the mapping  $T$  without trimming has the above optimal robust property.

**Remark 3.3.** The index  $n$  in the statement of Theorem 3.1 simply parameterizes how close  $Q \in B_H(P_0, r/\sqrt{n})$  and  $P_0$  are, and does not have to be interpreted as the sample size. The next theorem, however, is concerned with MSE's and the index  $n$  represents the sample size there.

The next theorem is our main result, which is concerned with (the supremum of) the MSE of the minimum Hellinger distance estimator  $\hat{\theta} = T(P_n)$  and other competing estimators. Let

$$(3.3) \quad \bar{B}_H(P_0, r/\sqrt{n}) = B_H(P_0, r/\sqrt{n}) \cap \left\{ Q \in \mathcal{M} : E_Q \left[ \sup_{\theta \in \Theta} |g(x, \theta)|^\eta \right] < \infty \right\}.$$

We use the notation  $P^{\otimes n}$  to denote the  $n$ -fold product measure of a probability measure  $P$ .

**Theorem 3.2.** Suppose that Assumption 3.1 holds.

(i): For every Fisher consistent and regular mapping  $T_a$ ,

$$\lim_{b \rightarrow \infty} \liminf_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n (\tau \circ T_a(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} \geq (1 + 4r^2) B^*,$$

for each  $r > 0$ .

(ii): The mapping  $T$  is Fisher consistent and regular, and the MHDE  $\hat{\theta} = T(P_n)$  satisfies

$$\lim_{b \rightarrow \infty} \lim_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n (\tau \circ T(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} = (1 + 4r^2) B^*,$$

for each  $r > 0$ .

**Remark 3.4.** This theorem establishes an asymptotic minimax optimality property of the MHDE in terms of MSE among all the estimators that satisfies the two conditions in Definition 3.1. Note that the expression  $\sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n (\tau \circ T_a(P_n) - \tau(\theta_0))^2 dQ^{\otimes n}$  is the maximum MSE of  $T_a(P_n)$  in a finite sample. Thus our criterion for evaluating  $T_a$  (and  $T$ ) is the limit of its maximum finite sample MSE. Taking the supremum over  $B_H$  before letting  $n$  go to infinity is important for capturing finite sample robustness properties. The method of calculating the truncated MSE first, then letting  $b \rightarrow \infty$ , is standard in the literature of robust estimation; see, for example, Bickel (1981). Once again, we are able to derive a sharp lower bound for the maximum MSE and show that it is achieved by the MHDE  $\hat{\theta} = T(P_n)$ .

**Remark 3.5.** Unlike in Theorem 3.1, optimality is achieved by the untrimmed version of the MHDE. Note that  $T(P_n)$  exists for large  $n$  under Assumption 3.1, in contrast to our discussion in Remark 3.2 on Theorem 3.1. Theorem 3.2, however, restricts the robustness neighborhood by an extra requirement as in (3.3). This is useful in showing that the untrimmed MHDE achieves the lower bound.

**Remark 3.6.** Theorem 3.2 proves that the MHDE is asymptotically optimally robust over a sequence of infinitesimal neighborhoods. Note that the Hellinger neighborhood over which the maximum of MSE is taken is nonparametric, in the sense that potential deviations from  $P_0$  cannot be indexed by a finite dimensional parameter. That is, our robustness concept demands uniform robustness over a nonparametric, infinitesimal neighborhood. The use of infinitesimal neighborhoods where the radius of the Hellinger ball shrinks at the rate  $n^{1/2}$  is useful in balancing the magnitude of bias and variance in our asymptotics. If one uses a fixed and global neighborhood, then the bias term would dominate the behavior of estimators. This may fail to provide a good approximation of finite sample behavior in actual applications, since in reality it would be reasonable to be concerned with both the stochastic fluctuation of estimators and their deterministic bias due to, say, data contamination. We note that there is a related but distinct literature on the asymptotics theory when the model is globally misspecified, as in White (1982), who considered parametric MLE. Kitamura (1998) and Kitamura (2002) offer such analysis for conditional and unconditional moment condition models. Moreover, Schennach (2007) provides novel and potentially very useful results of EL estimators and its variants in misspecified moment condition models. We regard our paper as a complement to, rather than a substitute for the results obtained in these papers. There are fundamental differences between the characteristics of the problems the current paper considers and those of the papers on misspecification. First, our object of interest is  $\theta_0$ , not a pseudo-true value, as we consider data perturbation rather than model misspecification. Second, the nature of our analysis is local and therefore the parameter value  $\theta_0$  in (1.1) is still identified asymptotically. Third, as noted above, we consider uniform robustness over a nonparametric neighborhood. The papers cited above consider pointwise problems. Therefore our approach deals with phenomena that are very different from the ones analyzed in the literature of misspecified models.

**Remark 3.7.** We have seen in Remark 2.1 that the Hellinger neighborhood  $B_H$  has nice and distinct properties, in particular the inclusion relationship (2.5). The neighborhood  $B_H$  is commonly used in the literature of robust estimation (of parametric models); see, for example, Beran (1977), Bickel (1981), and Rieder (1994). We should note, however, that other neighborhood systems have been

used in the literature as well. For example, one may replace the Hellinger distance  $H$  with the Kolmogorov-Smirnov (KS) distance in the definition of  $B_H$ . As Beran (1984) notes, however, that in order to guarantee robustness in the Kolmogorov-Smirnov neighborhood system one needs

“to use minimum distance estimates based on the Kolmogorov-Smirnov metric or a distance weaker than the Kolmogorov-Smirnov metric ... The general principle here is that the estimation distance be no stronger than the distance describing the contamination neighborhood...”

Donoho and Liu (1988) develop a general theory of the above point. What this means is that an estimator that is robust against perturbations within Kolmogorov-Smirnov neighborhoods has to be minimizing the KS (or weaker) distance. The “minimum KS estimator” for the moment restriction model would be indeed robust, but it cannot be semiparametrically efficient when the model assumption holds. Therefore, unlike the moment restriction MHDE, the estimator is not “robust and efficient.” Another drawback is its computation, since, unlike the moment restriction MHDE, no convenient algorithm to minimize the Kolmogorov-Smirnov distance under the moment restriction is known in the literature.

The above MSE theorem conveniently summarizes the desirable robustness properties of the MHDE in terms of both (deterministic) bias and variance. It has, however, some limitations. First, its minimaxity result is obtained within Fisher consistent and regular estimators. While these requirements are weak, it might be of interest to expand the class of estimators. More importantly, implicit in the MSE-based analysis is that we are interested in  $L^2$ -loss. One may wish to use other types of loss functions, however, and it is of interest to see whether the above minimax results can be extended to a larger class of loss. The next theorem addresses these two issues. Of course, the MSE has an advantage of subsuming the bias and the variance in one measure. To deal with general loss functions, the next theorem focuses on the risk of estimators around a Fisher-consistent mapping evaluated at the perturbed measure  $Q$ . This can be regarded as calculating the risk of the first bracket of the decomposition (1.2), that is, the stochastic part of the deviation of the estimator from the parameter of interest  $\theta_0$ .

Let  $\mathcal{S}$  be a set of all estimators, that is, the set of all  $\bar{\mathbb{R}}^p$ -valued measurable functions of the data  $(x_1, \dots, x_n)$ . We now investigate robust risk properties of this large class of estimators. The loss function we consider satisfies the following weak requirements.

**Assumption 3.2.** *The loss function  $\ell : \bar{\mathbb{R}}^p \rightarrow [0, \infty]$  is (i) symmetric subconvex (i.e., for all  $z \in \mathbb{R}^p$  and  $c \in \mathbb{R}$ ,  $\ell(z) = \ell(-z)$  and  $\{z \in \mathbb{R}^p : \ell(z) \leq c\}$  is convex); (ii) upper semicontinuous at infinity; and (iii) continuous on  $\bar{\mathbb{R}}^p$ .*

We now present an optimal risk property for the MHDE.

**Theorem 3.3.** Suppose that Assumptions 3.1 and 3.2 hold.

(i): For every Fisher consistent mapping  $T_a$ ,

$$\lim_{b \rightarrow \infty} \lim_{r \rightarrow \infty} \liminf_{n \rightarrow \infty} \inf_{S_n \in \mathcal{S}} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge \ell(\sqrt{n}(S_n - \tau \circ T_a(Q))) dQ^{\otimes n} \geq \int \ell dN(0, B^*).$$

(ii): The mapping  $T$  is Fisher consistent and the MHDE  $\hat{\theta} = T(P_n)$  satisfies

$$\lim_{b \rightarrow \infty} \lim_{r \rightarrow \infty} \lim_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge \ell(\sqrt{n}(\tau \circ T(P_n) - \tau \circ \bar{T}(Q))) dQ^{\otimes n} = \int \ell dN(0, B^*).$$

Theorem 3.3(ii) remains valid if  $T(P_n)$  is replaced by  $\bar{T}(P_n)$ . This theorem shows that the MHDE is once again optimally robust even for the general risk criterion, and this holds in the class of essentially all possible estimators. As noted above, the result is concerned with the stochastic component of the decomposition (1.2). Recall Theorem 3.1 has already established that the MHDE is optimal in terms of its bias, that is, the deterministic part of the decomposition (1.2) in the second bracket. The latter result does not depend on a specific loss function. Thus the MHDE enjoys general optimal robust properties under a quite general setting, both in terms of the stochastic component and the deterministic component. Note that analyzing these two parts separately is common in the literature of robust statistics: see, for example, Rieder (1994).

#### 4. SIMULATION

The purpose of this section is to examine the robustness properties of the MHDE and other well-known estimators such as GMM using Monte Carlo simulations. MATLAB is used for computation throughout the experiments. The sample size  $n$  is 100 for all designs, and we ran 5000 replications for each design.

**4.1. Experiment 1.** The baseline simulation design in this experiment follows that of Hall and Horowitz (1996). We then “contaminate” the simulated data to explore robustness of estimators. More specifically, let  $x = (x_1, x_2)' \sim N(0, 0.4^2 I_2)$ . This normal law corresponds to  $P_0$  in the

preceding sections. The specification of the moment function  $g$  is

$$g(x, \theta) = (\exp \{-0.72 - \theta(x_1 + x_2) + 3x_2\} - 1) \begin{pmatrix} 1 \\ x_2 \end{pmatrix}.$$

The moment condition  $\int g(x, \theta) dP_0 = 0$  is uniquely solved at  $\theta_0 = 3$ . The goal is to estimate this value using the above specification of  $g$  from contaminated data, which consists of 100 IID draws of  $x^* = (x_1^*, x_2^*)'$  generated according to

$$x_1^* = \begin{cases} x_1 & \text{with probability } 0.95, \\ x_1 + c \cdot \omega & \text{with probability } 0.05, \end{cases}$$

$$x_2^* = x_2$$

where  $\omega = \rho x_1 + \sqrt{1 - \rho^2} \cdot 0.4\xi$ . The contaminating variable  $\xi$  are specified to be either normal,  $\chi_1^2$ ,  $-\chi_1^2$  or  $t_3$ , though all of them are normalized to have mean zero and variance one. Note that  $\xi$  can be characterized to be a classical measurement error for  $\rho = 0$ , but not for the case with  $\rho = -0.5$ . We consider five estimators: empirical likelihood (EL), MHDE, exponential tilting (ET), GMM (GMM2) and CUE. GMM2 is calculated following the standard two step procedure where the initial estimate is obtained from identity weighting. The results are displayed in Table 1.<sup>2</sup>

While RMSE is a potentially informative measure, it can be a highly misleading as some of the estimators may not have finite moments. We thus focus on the simulated probability of an estimator deviating from the target  $\theta_0 = 3$  by more 0.5. The case with  $c = 0$  is the baseline without contamination. All the estimators work reasonably well for this case, though CUE seems to be somewhat problematic. As we add  $\xi$  with  $\rho = 0$ , interesting patterns emerge. Most notably, the performance of GMM2, which exhibits reasonable behavior with  $c = 0$ , deteriorates very rapidly as the DGP becomes perturbed. This casts serious doubt on the notion that GMM is a robust procedure. CUE is also sensitive to perturbations and performs poorly in general. In contrast, EL, MHDE and ET seem to be stable overall. Note, however, EL yields relatively high deviation probabilities in cases with  $c = 2$  and  $\xi \sim N(0, 1)$  or  $-\chi_1^2$ . A similar pattern appears with the case of negatively correlated errors (i.e.  $\rho = -0.5$ ), though in this case CUE is worse than GMM.

---

<sup>2</sup>When the values of moment function  $g$  fails to span the zero vector, EL, MHDE and ET cannot be calculated as they do not permit negative probability weights to set the moment condition at zero (see, for example, Kitamura (2006)). This occurs infrequently in our setting. Indeed, it never occurred in more than half of the simulation designs in our experiments and very few in others, except for a small number of designs with large  $c$  where the rate of its occurrence was at most about 1.5%. Those draws were discarded in calculating summary statistics in this section.



| c             | $\xi$       | RMSE  |       |       |       |       | $Pr \left\{ \left  \hat{\theta} - \theta_0 \right  > 0.5 \right\}$ |       |       |       |       |
|---------------|-------------|-------|-------|-------|-------|-------|--|-------|-------|-------|-------|
|               |             | EL    | MHDE  | ET    | GMM2  | CUE   | EL   | MHDE  | ET    | GMM2  | CUE   |
| $\rho = 0$    |             |       |       |       |       |       |  |       |       |       |       |
| 0             |             | 0.292 | 0.295 | 0.303 | 0.427 | 2.398 | 0.080  | 0.085 | 0.091 | 0.119 | 0.232 |
| 0.5           | $N$         | 0.304 | 0.304 | 0.311 | 0.442 | 2.265 | 0.095  | 0.097 | 0.099 | 0.140 | 0.222 |
| 1             | $N$         | 0.439 | 0.428 | 0.423 | 0.660 | 2.067 | 0.225  | 0.213 | 0.209 | 0.345 | 0.316 |
| 2             | $N$         | 0.777 | 0.704 | 0.678 | 1.291 | 2.066 | 0.451  | 0.375 | 0.357 | 0.667 | 0.493 |
| 0.5           | $\chi_1^2$  | 0.288 | 0.291 | 0.298 | 0.419 | 2.241 | 0.073  | 0.075 | 0.080 | 0.105 | 0.200 |
| 1             | $\chi_1^2$  | 0.295 | 0.295 | 0.297 | 0.383 | 2.045 | 0.086  | 0.087 | 0.089 | 0.119 | 0.185 |
| 2             | $\chi_1^2$  | 0.487 | 0.476 | 0.470 | 0.649 | 2.497 | 0.340  | 0.324 | 0.314 | 0.494 | 0.440 |
| 0.5           | $-\chi_1^2$ | 0.363 | 0.359 | 0.366 | 0.570 | 2.247 | 0.132  | 0.132 | 0.133 | 0.206 | 0.251 |
| 1             | $-\chi_1^2$ | 0.533 | 0.494 | 0.484 | 0.915 | 2.146 | 0.247  | 0.219 | 0.210 | 0.389 | 0.326 |
| 2             | $-\chi_1^2$ | 0.792 | 0.698 | 0.675 | 1.320 | 2.067 | 0.378  | 0.325 | 0.315 | 0.558 | 0.430 |
| 0.5           | $t_3$       | 0.318 | 0.317 | 0.324 | 0.486 | 2.291 | 0.102  | 0.101 | 0.105 | 0.151 | 0.225 |
| 1             | $t_3$       | 0.407 | 0.394 | 0.389 | 0.659 | 2.093 | 0.173  | 0.163 | 0.162 | 0.267 | 0.265 |
| 2             | $t_3$       | 0.658 | 0.603 | 0.582 | 1.100 | 2.078 | 0.346  | 0.310 | 0.297 | 0.529 | 0.407 |
| $\rho = -0.5$ |             |       |       |       |       |       |  |       |       |       |       |
| 1             | $N$         | 0.297 | 0.301 | 0.312 | 0.448 | 2.263 | 0.094  | 0.094 | 0.104 | 0.128 | 0.228 |
| 2             | $N$         | 0.320 | 0.320 | 0.325 | 0.446 | 2.095 | 0.116  | 0.118 | 0.122 | 0.152 | 0.234 |
| 1             | $\chi_1^2$  | 0.297 | 0.303 | 0.314 | 0.452 | 2.443 | 0.084  | 0.096 | 0.110 | 0.136 | 0.250 |
| 2             | $\chi_1^2$  | 0.286 | 0.291 | 0.299 | 0.427 | 2.313 | 0.084  | 0.086 | 0.090 | 0.102 | 0.212 |
| 1             | $-\chi_1^2$ | 0.311 | 0.313 | 0.319 | 0.498 | 2.525 | 0.104  | 0.104 | 0.112 | 0.158 | 0.262 |
| 2             | $-\chi_1^2$ | 0.404 | 0.391 | 0.386 | 0.640 | 2.436 | 0.160  | 0.150 | 0.148 | 0.232 | 0.293 |
| 1             | $t_3$       | 0.298 | 0.299 | 0.306 | 0.474 | 2.480 | 0.076  | 0.078 | 0.088 | 0.114 | 0.232 |
| 2             | $t_3$       | 0.340 | 0.335 | 0.339 | 0.556 | 2.124 | 0.104  | 0.106 | 0.110 | 0.162 | 0.234 |

TABLE 1. The second column “ $\xi$ ” specifies the distribution of  $\xi$ , where the labels  $N$ ,  $\chi_1^2$ ,  $-\chi_1^2$ , and  $t_3$  denote  $N(0, 1)$ ,  $(\chi_1^2 - 1)/\sqrt{2}$ ,  $-(\chi_1^2 - 1)/\sqrt{2}$ , and Student- $t_3/\sqrt{3}$ , respectively.

**4.2. Experiment 2.** This experiment uses the same model specification of  $g$  as above, though the DGP is replaced by a family of normal distributions. Experiment 1 employs two types of perturbations ( $\rho = 0$  and  $\rho = -0.5$ ) with varied magnitudes controlled by the parameter  $c$ , whereas this experimental

setting attempts to perturb the original DGP  $x \sim N(0, 0.4^2 I_2)$  into different directions. More specifically, we use  $x \sim N(0, \Sigma_{(\delta, \rho)})$ , where

$$\Sigma_{(\delta, \rho)} = 0.4^2 \begin{pmatrix} (1 + \delta)^2 & \rho(1 + \delta) \\ \rho(1 + \delta) & 1 \end{pmatrix}.$$

The unperturbed case thus corresponds to  $\delta = \rho = 0$ . In the simulation we set  $\rho = 0.1\sqrt{2} \cos(2\pi\omega)$  and  $\delta = 0.1 \sin(2\pi\omega)$  and let  $\omega$  vary over  $\omega_j = j/64, j = 0, \dots, 63$ . This yields 64 different designs, for each of them 5000 replications is performed and RMSE and  $Pr\left\{\left|\hat{\theta} - \theta_0\right| > 0.5\right\}$  are calculated. The results are presented in Figure 4.1. In the upper panel, each curve represents the RMSE of a particular estimator as a function of  $\omega_j$ . The lower panel (labeled “Pr”) displays deviation probabilities.

We again focus on the results for deviation probabilities due to our concern about the existence of moments. The graph labeled “Pr(all)” shows that the CUE’s performance is extremely sensitive to data perturbations considered here; its deviation probability is uniformly higher than those of the rest, and the difference can be quite large in places. We therefore plotted the same results without CUE on the graph labeled “Pr(no CUE)” to visualize the relative rankings of the other four estimators more clearly. We see that GMM2 is affected by perturbations much more than EL, MHDE and ET except for the values of  $\omega$ ’s between 0.4 and 0.6, where the performance of the four estimators are rather close. ET seems to perform a little worse than MHDE and EL.

One needs be cautious in drawing conclusions based on limited simulation experiments as presented here. Nevertheless, it appears that two general features emerge from our results. First, the GMM type estimators (two step GMM and CUE) tend to be highly sensitive to data perturbations. Applying Beran’s (1977) logic that connects the robustness of estimators to the forms of their objective functions, this may be attributed to the fact the GMM objective function is quadratic and therefore tends to react sensitively to the added noises. Second, EL, MHDE and ET are relatively well-behaved, and their rankings, not surprisingly, vary depending on the simulation design. The performance of MHDE, however, seems more stable compared with that of EL or ET: EL and ET exhibits more instability in Experiment 1 and Experiment 2, respectively. Note that EL, MHDE, ET and CUE correspond to the GEL estimator with  $\gamma = -1, -\frac{1}{2}, 0, 1$  in equation (2.6) of Newey and Smith (2004). Given the good theoretical robustness property of the MHDE and the proximity of EL and ET in terms of their  $\gamma$  values, it is interesting to observe the reasonably robust behavior of EL and ET. Note that CUE, whose behavior is quite different from that of the MHDE and thus highly non-robust, has  $\gamma = 1$ , a value that is much higher than the optimally robust  $\gamma = -1/2$  of the MHDE.

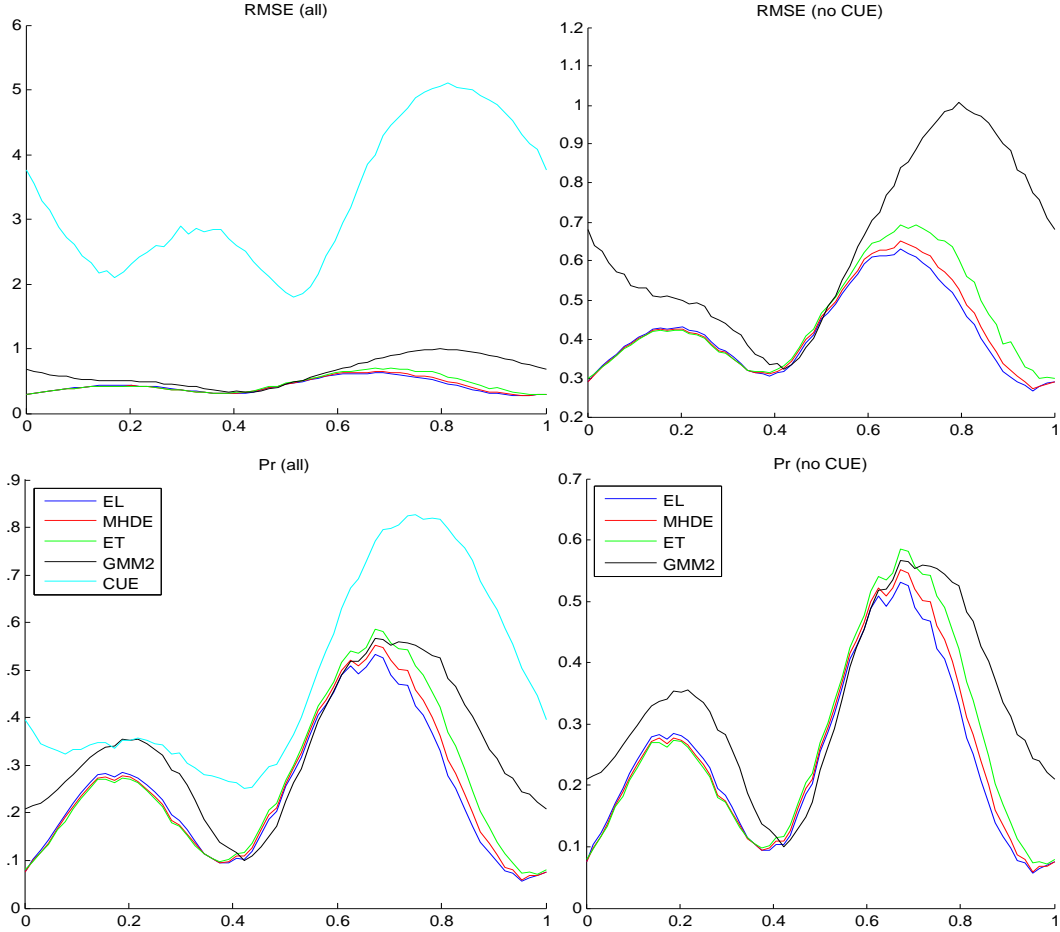


FIGURE 4.1. Local Neighborhood of the True Model. “Pr” denotes  $Pr \left\{ \left| \hat{\theta} - \theta_0 \right| > 0.5 \right\}$ .

## 5. CONCLUSION

In this paper we have explored the issue of robust estimation in a moment restriction model. The model is semiparametric and distribution-free, therefore imposes mild assumptions. Yet it is reasonable to expect that the probability law of observations may have some deviations from the ideal distribution being modeled. It is then sensible to seek estimation procedures that are robust against

slight perturbations in the probability measure that generates observations, which can be caused by, for example, data contamination. Our main theoretical result is that the MHDE possesses optimal minimax robust properties. We show this by deriving three asymptotic minimax theorems concerning bias, MSE and general risk, and in each criterion the MHDE achieves the minimax lower bound. Moreover, it remains semiparametrically efficient when the model assumptions hold. Convenient numerical algorithms for its implementation are provided. Our simulation results indicate that GMM can be highly sensitive to data perturbations. The performance of the MHDE remains stable over a wide range of simulation designs, which is in accordance with our theoretical findings.

The results obtained in this paper are concerned with estimation, though it might be possible to extend our robustness theory to parameter testing problems. Interestingly, there exists a literature on parametric robust inference based on the MHDE method. It is of practical importance to consider robust methods for parameter testing and confidence interval calculations so that the results of statistical inference for moment restriction models are reliable and not too sensitive to departures from model assumptions. We plan to investigate robust testing procedure in moment condition models in our future research.

## 6. APPENDIX

This Appendix presents the proofs of some of the results presented in the previous sections.

**Notation.** Let  $C$  be a generic positive constant, and  $\|\cdot\|$  be the  $L_2$ -metric,

$$\begin{aligned} \theta_n &= \theta_0 + t/\sqrt{n}, \quad \bar{T}_{Q_n} = \bar{T}(Q_n), \quad \bar{T}_{P_n} = \bar{T}(P_n), \\ \bar{P}_{\theta,Q} &= \arg \min_{P \in \bar{\mathcal{P}}_\theta} H(P, Q), \quad R_n(Q, \theta, \gamma) = - \int \frac{1}{(1 + \gamma' g_n(x, \theta))} dQ, \\ g_n(x, \theta) &= g(x, \theta) \mathbb{I}\{x \in \mathcal{X}_n\}, \quad \Lambda_n = G' \Omega^{-1} g_n(x, \theta_0), \quad \Lambda = G' \Omega^{-1} g(x, \theta_0), \\ \psi_{n, Q_n} &= -2 \left( \int \Lambda_n \Lambda_n' dQ_n \right)^{-1} \int \Lambda_n \left\{ dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} \right\} dQ_n^{1/2}. \end{aligned}$$

### 6.1. Proof of Theorem 3.1.

6.1.1. *Proof of (i).* Pick arbitrary  $r > 0$  and  $t \in \mathbb{R}^p$ . Consider the following parametric submodel having the likelihood ratio

$$(6.1) \quad \frac{dP_{\theta_n, \zeta_n}}{dP_0} = \frac{1 + \zeta_n' g_n(x, \theta_n)}{\int (1 + \zeta_n' g_n(x, \theta_n)) dP_0} = f(x, \theta_n, \zeta_n),$$

where

$$\zeta_n = -E_{P_0} [g(x, \theta_n) g_n(x, \theta_n)']^{-1} E_{P_0} [g(x, \theta_n)].$$

Note that  $P_{\theta_0,0} = P_0$ ,  $P_{\theta_n,\zeta_n} \in \mathcal{P}_{\theta_n}$  (by the definition of  $\zeta_n$ ), and  $\zeta_n = O(n^{-1/2})$  (by the proof of Lemma 7.4 (i)). Also, since  $\sup_{x \in \mathcal{X}} |\zeta'_n g_n(x, \theta_n)| = O(n^{-1/2} m_n) = o(1)$ , the likelihood ratio  $\frac{dP_{\theta_n,\zeta_n}}{dP_0}$  is well-defined for all  $n$  large enough. So, for this submodel the mapping  $T_a$  must satisfy (3.1).

We now evaluate the Hellinger distance between  $P_{\theta_n,\zeta_n}$  and  $P_0$ . An expansion around  $\zeta_n = 0$  yields

$$H(P_{\theta_n,\zeta_n}, P_0) = \left\| \zeta'_n \frac{\partial f(x, \theta_n, \zeta_n)^{1/2}}{\partial \zeta_n} \Big|_{\zeta_n=0} dP_0^{1/2} + \frac{1}{2} \zeta'_n \frac{\partial^2 f(x, \theta_n, \zeta_n)^{1/2}}{\partial \zeta_n \partial \zeta'_n} \Big|_{\zeta_n=\dot{\zeta}_n} \zeta_n dP_0^{1/2} \right\|,$$

where  $\dot{\zeta}_n$  is a point on the line joining  $\zeta_n$  and 0, and

$$\frac{\partial f(x, \theta_n, \zeta_n)^{1/2}}{\partial \zeta_n} \Big|_{\zeta_n=0} = \frac{1}{2} \{g_n(x, \theta_n) - E_{P_0}[g_n(x, \theta_n)]\},$$

$$\begin{aligned} \frac{\partial^2 f(x, \theta_n, \zeta_n)^{1/2}}{\partial \zeta_n \partial \zeta'_n} &= -\frac{1}{4} (1 + \zeta'_n g_n(x, \theta_n))^{-3/2} (1 + \zeta'_n E_{P_0}[g_n(x, \theta_n)])^{-1/2} g_n(x, \theta_n) g_n(x, \theta_n)' \\ &\quad - \frac{1}{2} (1 + \zeta'_n g_n(x, \theta_n))^{-1/2} (1 + \zeta'_n E_{P_0}[g_n(x, \theta_n)])^{-3/2} g_n(x, \theta_n) E_{P_0}[g_n(x, \theta_n)]' \\ &\quad + \frac{3}{4} (1 + \zeta'_n g_n(x, \theta_n))^{1/2} (1 + \zeta'_n E_{P_0}[g_n(x, \theta_n)])^{-5/2} E_{P_0}[g_n(x, \theta_n)] E_{P_0}[g_n(x, \theta_n)]'. \end{aligned}$$

Thus, a lengthy but straightforward calculation combined with Lemma 7.4,  $\zeta_n = O(n^{-1/2})$ , and  $\sup_{x \in \mathcal{X}} |\zeta'_n g_n(x, \theta_n)| = o(1)$  implies

$$(6.2) \quad nH(P_{\theta_n,\zeta_n}, P_0)^2 = n \left\| \frac{1}{2} \zeta'_n (g_n(x, \theta_n) - E_{P_0}[g_n(x, \theta_n)]) dP_0^{1/2} \right\|^2 + o(1) \rightarrow \frac{1}{4} t' \Sigma^{-1} t.$$

Based on this limit, a lower bound of the maximum bias of  $T_a$  is obtained as (see, Rieder (1994, eq. (56) on p. 180))

$$\begin{aligned} &\liminf_{n \rightarrow \infty} \sup_{Q \in B_H(P_0, r/\sqrt{n})} n (\tau \circ T_a(Q) - \tau(\theta_0))^2 \\ &\geq \liminf_{n \rightarrow \infty} \sup_{\{t \in \mathbb{R}^p : P_{\theta_n,\zeta_n} \in B_H(P_0, r/\sqrt{n})\}} n (\tau \circ T_a(P_{\theta_n,\zeta_n}) - \tau(\theta_0))^2 \\ &\geq \max_{\{t \in \mathbb{R}^p : \frac{1}{4} t' \Sigma t \leq r^2 - \epsilon\}} \left( \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' t \right)^2 = 4(r^2 - \epsilon) \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' \Sigma^{-1} \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right), \end{aligned}$$

for each  $\epsilon \in (0, r^2)$ , where the first inequality follows from the set inclusion relationship, the second inequality follows from (3.1) and (6.2), and the equality follows from the Kuhn-Tucker theorem. Since  $\epsilon$  can be arbitrarily small, we obtain the conclusion.

6.1.2. *Proof of (ii).* Pick arbitrary  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ . We first show the Fisher consistency of  $\bar{T}$ . From Lemma 7.2 (note:  $P_{\theta_n, \zeta_n} \in B_H(P_0, r/\sqrt{n})$  for all  $n$  large enough),

$$\begin{aligned} \sqrt{n}(\bar{T}(P_{\theta_n, \zeta_n}) - \theta_0) &= -\sqrt{n}\Sigma^{-1} \int \Lambda_n dP_{\theta_n, \zeta_n} + o(1) \\ &= \Sigma^{-1} G' \Omega^{-1} \int \partial g(x, \dot{\theta}) / \partial \theta dP_{\theta_n, \zeta_n} t + o(1) \\ &\rightarrow t \end{aligned}$$

for all  $n$  large enough, where  $\dot{\theta}$  is a point on the line joining  $\theta_n$  and  $\theta_0$ , the second equality follows from  $\int g(x, \theta_0) \mathbb{I}\{x \notin \mathcal{X}_n\} dP_{\theta_n, \zeta_n} = o(n^{-1/2})$  (by a similar argument to (7.2)),  $\int g(x, \theta_n) dP_{\theta_n, \zeta_n} = 0$  (by  $P_{\theta_n, \zeta_n} \in \mathcal{P}_{\theta_n}$ ), and an expansion around  $\theta_n = \theta_0$ , and the convergence follows from the last statement of Lemma 7.4 (i). Therefore,  $\bar{T}$  is Fisher consistent.

We next show (3.1). An expansion of  $\tau \circ \bar{T}_{Q_n}$  around  $\bar{T}_{Q_n} = \theta_0$ , Lemmas 7.1 (ii) and 7.2, and Assumption 3.1 (viii) imply

$$\begin{aligned} \sqrt{n}(\tau \circ \bar{T}_{Q_n} - \tau(\theta_0)) &= -\sqrt{n} \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' \Sigma^{-1} \int \Lambda_n dQ_n + o(1) \\ &= -\sqrt{n} \nu'_0 \int \Lambda_n \{dQ_n^{1/2} - dP_0^{1/2}\} dQ_n^{1/2} - \sqrt{n} \nu'_0 \int \Lambda_n dP_0^{1/2} \{dQ_n^{1/2} - dP_0^{1/2}\} + o(1), \end{aligned}$$

where we denote  $\nu'_0 = \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' \Sigma^{-1}$ . From the triangle inequality,

$$\begin{aligned} &n(\tau \circ \bar{T}_{Q_n} - \tau(\theta_0))^2 \\ &\leq n \left\{ \left| \nu'_0 \int \Lambda_n \{dQ_n^{1/2} - dP_0^{1/2}\} dQ_n^{1/2} \right|^2 + \left| \nu'_0 \int \Lambda_n \{dQ_n^{1/2} - dP_0^{1/2}\} dP_0^{1/2} \right|^2 \right. \\ &\quad \left. + 2 \left| \nu'_0 \int \Lambda_n \{dQ_n^{1/2} - dP_0^{1/2}\} dQ_n^{1/2} \right| \left| \nu'_0 \int \Lambda_n \{dQ_n^{1/2} - dP_0^{1/2}\} dP_0^{1/2} \right| \right\} + o(1) \\ &= n \{A_1 + A_2 + 2A_3\} + o(1). \end{aligned}$$

For  $A_1$ , observe that

$$A_1 \leq \left| \nu'_0 \int \Lambda_n \Lambda'_n dQ_n \nu_0 \right| \left| \int \{dQ_n^{1/2} - dP_0^{1/2}\}^2 \right| \leq B^* \frac{r^2}{n} + o(n^{-1}),$$

where the first inequality follows from the Cauchy-Schwarz inequality, and the second inequality follows from Lemma 7.5 (i) and  $Q_n \in B_H(P_0, r/\sqrt{n})$ . Similarly, we have  $A_2 \leq B^* \frac{r^2}{n} + o(n^{-1})$  and  $A_3 \leq B^* \frac{r^2}{n} + o(n^{-1})$ . Combining these terms,

$$\limsup_{n \rightarrow \infty} n(\tau \circ \bar{T}_{Q_n} - \tau(\theta_0))^2 \leq 4r^2 B^*,$$

for any sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$  and  $r > 0$ . This implies the conclusion because  $B_H(P_0, r/\sqrt{n})$  is compact with respect to the Hellinger distance for each  $n \in \mathbb{N}$  and  $\tau \circ \bar{T}(Q)$  is upper semi-continuous at each  $Q \in \mathcal{M}$  under the Hellinger distance (Lemma 7.1 (i)).

## 6.2. Proof of Theorem 3.2.

6.2.1. *Proof of (i).* Pick arbitrary  $\epsilon \in (0, r^2)$  and  $r > 0$ . Consider the parametric submodel  $P_{\theta_n, \zeta_n}$  defined in (6.1). The convolution theorem (Theorem 25.20 of van der Vaart (1998)) implies that for each  $t \in \mathbb{R}^p$ , there exists a probability measure  $M_0$  which does not depend on  $t$  and satisfies

$$(6.3) \quad \sqrt{n}(\tau \circ T_a(P_n) - \tau \circ T_a(P_{\theta_n, \zeta_n})) \xrightarrow{d} M_0 * N(0, B^*) \quad \text{under } P_{\theta_n, \zeta_n},$$

where the symbol  $*$  denote convolution. Let

$$t^* = \arg \max_{\{t \in \mathbb{R}^p: \frac{1}{4}t' \Sigma t \leq r^2 - \epsilon\}} \left( \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' t \right)^2 \quad \text{s.t.} \quad \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' t^* \int \xi dM_0 * N(0, B^*) \geq 0.$$

Since the integral  $\int \xi dM_0 * N(0, B^*)$  does not depend on  $t$ , such  $t^*$  always exists. From  $\frac{1}{4}t^{*'} \Sigma t^* \leq r^2 - \epsilon$  and (6.2), it holds that  $P_{\theta_0 + t^*/\sqrt{n}, \zeta_n} \in B_H(P_0, r/\sqrt{n})$  for all  $n$  large enough. Also, note that  $E_{P_{\theta_n, \zeta_n}}[\sup_{\theta \in \Theta} |g(x, \theta)|^q] < \infty$  for all  $n$  large enough (by  $\sup_{x \in \mathcal{X}} |\zeta_n' g_n(x, \theta_n)| = o(1)$  and Assumption 3.1 (v)). Thus,  $P_{\theta_0 + t^*/\sqrt{n}, \zeta_n} \in \bar{B}_H(P_0, r/\sqrt{n})$  for all  $n$  large enough, and we have

$$\begin{aligned} & \lim_{b \rightarrow \infty} \liminf_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n (\tau \circ T_a(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} \\ & \geq \lim_{b \rightarrow \infty} \liminf_{n \rightarrow \infty} \int b \wedge n (\tau \circ T_a(P_n) - \tau(\theta_0))^2 dP_{\theta_0 + t^*/\sqrt{n}, \zeta_n}^{\otimes n} \\ & = \lim_{b \rightarrow \infty} \liminf_{n \rightarrow \infty} \int b \wedge n \left( \xi + \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' t^* \right)^2 dM_0 * N(0, B^*) \\ & = \int \xi^2 dM_0 * N(0, B^*) + \left( \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' t^* \right)^2 + 2 \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' t^* \int \xi dM_0 * N(0, B^*) \\ & \geq \{1 + 4(r^2 - \epsilon)\} B^*, \end{aligned}$$

where the first equality follows from the Fisher consistency of  $T_a$ , (6.5), and the continuous mapping theorem, the second equality follows from the monotone convergence theorem, and the second inequality follows from the definition of  $t^*$ . Since  $\epsilon$  can be arbitrarily small, we obtain the conclusion.

6.2.2. *Proof of (ii).* Pick arbitrary  $r > 0$  and  $b > 0$ . Applying the inequality  $b \wedge (c_1 + c_2) \leq b \wedge c_1 + b \wedge c_2$  for any  $c_1, c_2 \geq 0$ ,

$$\begin{aligned}
& \limsup_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n (\tau \circ T(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} \\
\leq & \limsup_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n (\tau \circ T(P_n) - \tau \circ \bar{T}(P_n))^2 dQ^{\otimes n} \\
& + 2 \limsup_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge \{n |\tau \circ T(P_n) - \tau \circ \bar{T}(P_n)| |\tau \circ \bar{T}(P_n) - \tau(\theta_0)|\} dQ^{\otimes n} \\
& + \limsup_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n (\tau \circ \bar{T}(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} \\
(6.4) = & A_1 + 2A_2 + A_3,
\end{aligned}$$

For  $A_1$ ,

$$\begin{aligned}
A_1 & \leq b \times \limsup_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int_{(x_1, \dots, x_n) \notin \mathcal{X}_n^n} dQ^{\otimes n} \\
& \leq b \times \limsup_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \sum_{i=1}^n \int_{x_i \notin \mathcal{X}_n} dQ \\
(6.5) \quad & \leq b \times \limsup_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} nm_n^{-\eta} E_Q \left[ \sup_{\theta \in \Theta} |g(x, \theta)|^\eta \right] = 0,
\end{aligned}$$

where the first inequality follows from  $T(P_n) = \bar{T}(P_n)$  for all  $(x_1, \dots, x_n) \in \mathcal{X}_n^n$ , the second inequality follows from a set inclusion relation, the third inequality follows from the Markov inequality, and the equality follows from Assumption 3.1 (vii) and  $E_Q [\sup_{\theta \in \Theta} |g(x, \theta)|^\eta] < \infty$  for all  $Q \in \bar{B}_H(P_0, r/\sqrt{n})$ . Similarly, we have  $A_2 = 0$ .

We now consider  $A_3$ . Note that the mapping  $f_{b,n}(Q) = \int b \wedge n (\tau \circ \bar{T}(P_n) - \tau(\theta_0))^2 dQ^{\otimes n}$  is continuous in  $Q \in B_H(P_0, r/\sqrt{n})$  under the Hellinger distance for each  $n$ , and the set  $B_H(P_0, r/\sqrt{n})$  (not  $\bar{B}_H(P_0, r/\sqrt{n})$ ) is compact under the Hellinger distance for each  $n$ . Thus, there exists  $\tilde{Q}_{b,n} \in$



$B_H(P_0, r/\sqrt{n})$  such that  $\sup_{Q \in B_H(P_0, r/\sqrt{n})} f_n(Q) = f_n(\tilde{Q}_{b,n})$  for each  $n$ . Then we have

$$\begin{aligned}
A_3 &\leq \limsup_{n \rightarrow \infty} \sup_{Q \in B_H(P_0, r/\sqrt{n})} \int b \wedge n (\tau \circ \bar{T}(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} \\
&= \limsup_{n \rightarrow \infty} \int b \wedge n (\tau \circ \bar{T}(P_n) - \tau(\theta_0))^2 d\tilde{Q}_{b,n}^{\otimes n} \\
&= \int b \wedge (\xi + \tilde{t}_b)^2 dN(0, B^*) \\
&\leq B^* + \tilde{t}_b^2 \\
&\leq (1 + 4r^2) B^*,
\end{aligned}$$

where  $\tilde{t}_b = \limsup_{n \rightarrow \infty} \sqrt{n} (\tau \circ \bar{T}(\tilde{Q}_{b,n}) - \tau(\theta_0))$ , the first inequality follows from  $\bar{B}_H(P_0, r/\sqrt{n}) \subseteq B_H(P_0, r/\sqrt{n})$ , the second equality follows from Lemma 7.8 (with  $Q_n = \tilde{Q}_{b,n}$ ) and the continuous mapping theorem, the second inequality follows from  $b \wedge c \leq c$  and a direct calculation, and the last inequality follows from Theorem 3.1 (ii). Combining these results, the conclusion is obtained.

### 6.3. Proof of Theorem 3.3.

6.3.1. *Proof of (i).* Consider the parametric submodel  $P_{\theta_n, \zeta_n}$  defined in (6.1). Since  $\ell$  is uniformly continuous on  $\bar{\mathbb{R}}^p$  (by Assumption 3.2) and  $T_a$  is Fisher consistent,

$$b \wedge \ell(\sqrt{n} \{S_n - \tau \circ T_a(P_{\theta_n, \zeta_n})\}) - b \wedge \ell\left(\sqrt{n} \{S_n - \tau(\theta_0)\} - \left(\frac{\partial \tau(\theta_0)}{\partial \theta}\right)' t\right) \rightarrow 0,$$

uniformly in  $t$ ,  $|t| < c$  and  $\{S_n\}_{n \in \mathbb{N}}$  for each  $c > 0$  and  $b > 0$ . Thus,

(6.6)

$$\inf_{S_n \in \mathcal{S}} \sup_{|t| \leq c} \int b \wedge \ell(\sqrt{n} \{S_n - \tau \circ T_a(P_{\theta_n, \zeta_n})\}) dP_{\theta_n, \zeta_n}^{\otimes n} = \inf_{R_n \in \mathcal{R}} \sup_{|t| \leq c} \int b \wedge \ell\left(R_n - \left(\frac{\partial \tau(\theta_0)}{\partial \theta}\right)' t\right) dP_{\theta_n, \zeta_n}^{\otimes n} + o(1),$$

for each  $c > 0$ , where  $R_n = \sqrt{n} \{S_n - \tau(\theta_0)\}$  is a standardized estimator and  $\mathcal{R} = \{\sqrt{n} \{S_n - \tau(\theta_0)\} : S_n \in \mathcal{S}\}$ .

By expanding the log likelihood ratio  $\log \frac{dP_{\theta_n, \zeta_n}^{\otimes n}}{dP_0^{\otimes n}}$  around  $\zeta_n = 0$ ,

$$\begin{aligned}
\log \frac{dP_{\theta_n, \zeta_n}^{\otimes n}}{dP_0^{\otimes n}} &= \zeta_n' \sum_{i=1}^n \{g_n(x_i, \theta_n) - E_{P_0}[g_n(x, \theta_n)]\} \\
&\quad - \frac{\zeta_n' \sum_{i=1}^n g_n(x_i, \theta_n) g_n(x_i, \theta_n) \zeta_n}{2 \left(1 + \zeta_n' g_n(x_i, \theta_n)\right)^2} + \frac{n \zeta_n' E_{P_0}[g_n(x, \theta_n)] E_{P_0}[g_n(x, \theta_n)]' \zeta_n}{2 \left(1 + \zeta_n' \int g_n(x, \theta_n)\right)^2} \\
&= L_1 - L_2 + L_3.
\end{aligned}$$

where  $\check{\zeta}_n$  and  $\ddot{\zeta}_n$  are points on the line joining  $\zeta_n$  and 0. For  $L_1$ , an expansion of  $g_n(x, \theta_n)$  (in  $\zeta_n$ ) around  $\theta_n = \theta_0$  combined with Lemma 7.4 (i) implies that under  $P_0$ ,

$$L_1 = -t'G'\Omega^{-1}\frac{1}{\sqrt{n}}\sum_{i=1}^n\{g_n(x_i, \theta_n) - E_{P_0}[g_n(x, \theta_n)]\} + o_p(1).$$

Also, Lemma 7.4 (i) and  $\sup_{x \in \mathcal{X}} |\zeta_n' g_n(x, \theta_n)| = o(1)$  imply that under  $P_0$ ,

$$L_2 \xrightarrow{p} \frac{1}{2}t'\Sigma t, \quad L_3 \rightarrow 0.$$

Therefore, in the terminology of Rieder (1994, Definition 2.2.9), the parametric model  $P_{\theta_n, \zeta_n}$  is asymptotically normal with the asymptotic sufficient statistic  $-G'\Omega^{-1}\frac{1}{\sqrt{n}}\sum_{i=1}^n\{g_n(x_i, \theta_n) - E_{P_0}[g_n(x, \theta_n)]\}$  and the asymptotic covariance matrix  $\Sigma$ . Note that this is essentially the LAN (local asymptotic normality) condition introduced by LeCam. If  $P_{\theta_n, \zeta_n}$  is asymptotically normal in this sense, we can directly apply the result of the minimax risk bound by Rieder (1994, Theorem 3.3.8 (a)), that is

$$(6.7) \quad \lim_{b \rightarrow \infty} \lim_{c \rightarrow \infty} \liminf_{n \rightarrow \infty} \inf_{S_n \in \mathcal{S}} \sup_{|t| \leq c} \int b \wedge \ell \left( R_n - \left( \frac{\partial \tau(\theta_0)}{\partial \theta} \right)' t \right) dP_{\theta_n, \zeta_n}^{\otimes n} \geq \int \ell dN(0, B^*).$$

From (6.6) and (6.7),

$$\lim_{b \rightarrow \infty} \lim_{c \rightarrow \infty} \liminf_{n \rightarrow \infty} \inf_{S_n \in \mathcal{S}} \sup_{|t| \leq c} \int b \wedge \ell(\sqrt{n}\{S_n - \tau \circ T_a(P_{\theta_n, \zeta_n})\}) dP_{\theta_n, \zeta_n}^{\otimes n} \geq \int \ell dN(0, B^*).$$

Finally, since  $E_{P_{\theta_n, \zeta_n}}[\sup_{\theta \in \Theta} |g(x, \theta)|^\eta] < \infty$  for all  $n$  large enough (by  $\sup_{x \in \mathcal{X}} |\zeta_n' g_n(x, \theta_n)| = o(1)$  and Assumption 3.1 (v)), we have  $P_{\theta_n, \zeta_n} \in \bar{B}_H(P_0, r/\sqrt{n})$  for all  $t$  satisfying  $\frac{1}{4}t'\Sigma t \leq r^2 - \epsilon$  with any  $\epsilon \in (0, r^2)$  and all  $n$  large enough. Therefore, the set inclusion relation yields

$$\begin{aligned} & \lim_{b \rightarrow \infty} \lim_{r \rightarrow \infty} \liminf_{n \rightarrow \infty} \inf_{S_n \in \mathcal{S}} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge \ell(\sqrt{n}\{S_n - \tau \circ T_a(Q)\}) dQ^{\otimes n} \\ & \geq \lim_{b \rightarrow \infty} \lim_{c \rightarrow \infty} \liminf_{n \rightarrow \infty} \inf_{S_n \in \mathcal{S}} \sup_{|t| \leq c} \int b \wedge \ell(\sqrt{n}\{S_n - \tau \circ T_a(P_{\theta_n, \zeta_n})\}) dP_{\theta_n, \zeta_n}^{\otimes n}, \end{aligned}$$

which implies the conclusion.

6.3.2. *Proof of (ii).* Pick arbitrary  $r > 0$  and  $b > 0$ . Since  $T(P_n) = \bar{T}(P_n)$  for all  $(x_1, \dots, x_n) \in \mathcal{X}_n^n$ ,

$$(6.8) \quad \begin{aligned} & \lim_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge \ell(\sqrt{n}\{\tau \circ T(P_n) - \tau \circ \bar{T}(Q)\}) dQ^{\otimes n} \\ & \leq \lim_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int_{(x_1, \dots, x_n) \notin \mathcal{X}_n^n} b \wedge \ell(\sqrt{n}\{\tau \circ T(P_n) - \tau \circ \bar{T}(Q)\}) dQ^{\otimes n} \\ & \quad + \lim_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int_{(x_1, \dots, x_n) \in \mathcal{X}_n^n} b \wedge \ell(\sqrt{n}\{\tau \circ \bar{T}(P_n) - \tau \circ \bar{T}(Q)\}) dQ^{\otimes n}. \end{aligned}$$

An argument similar to (6.5) implies that the first term of (6.8) is zero. From  $\mathcal{X}_n^n \subseteq \mathcal{X}^n$  and  $\bar{B}_H(P_0, r/\sqrt{n}) \subseteq B_H(P_0, r/\sqrt{n})$ , the second term of (6.8) is bounded from above by

$$\lim_{n \rightarrow \infty} \sup_{Q \in B_H(P_0, r/\sqrt{n})} \int b \wedge \ell(\sqrt{n} \{\tau \circ \bar{T}(P_n) - \tau \circ \bar{T}(Q)\}) dQ^{\otimes n} = \int b \wedge \ell dN(0, B^*),$$

where the equality follows from Lemma 7.8, the uniform continuity of  $\ell$  over  $\bar{\mathbb{R}}^p$ , and compactness of  $B_H(P_0, r/\sqrt{n})$  under the Hellinger distance. Let  $r \rightarrow \infty$  and  $b \rightarrow \infty$  then the conclusion follows.

## 7. AUXILIARY LEMMAS

**Lemma 7.1.** *Suppose that Assumption 3.1 holds. Then*

- (i): *for each  $n \in \mathbb{N}$ ,  $\bar{T}(Q)$  exists and is upper semi-continuous at each  $Q \in \mathcal{M}$  under the Hellinger distance,*
- (ii):  *$\bar{T}_{Q_n} \rightarrow \theta_0$  as  $n \rightarrow \infty$  for each  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ .*

**Proof of (i).** The proof is based on Lemma 1 of Kitamura (2001). See also Beran (1984, p. 744). Pick an arbitrary  $n \in \mathbb{N}$ . Denote  $R_n(Q, \theta) = \inf_{P \in \bar{\mathcal{P}}_\theta} H(P, Q)$ . Since  $g_n(x, \theta)$  is bounded a.s. for all  $\theta \in \Theta$ , the duality of partially finite programming (Borwein and Lewis (1993)) implies that  $R_n(Q, \theta) = \max_{\gamma \in \mathbb{R}^m} R_n(Q, \theta, \gamma)$  for each  $(Q, \theta) \in \mathcal{M} \times \Theta$ . From Theorem 10.8 of Rockafeller (1970) and Assumption 3.1 (iv),  $R_n(Q, \theta)$  is continuous in  $(Q, \theta) \in \mathcal{M} \times \Theta$  under the Lévy metric (for  $\mathcal{M}$ ). This continuity also implies that for each  $Q \in \mathcal{M}$ ,  $R_n(Q, \theta)$  is continuous in  $\theta \in \Theta$ . Since  $\Theta$  is compact (Assumption 3.1 (ii)), the Weierstrass theorem guarantees the existence of  $\bar{T}(Q) = \arg \min_{\theta \in \Theta} R_n(Q, \theta)$ . Also, since  $R_n(Q, \theta)$  is continuous in  $(Q, \theta) \in \mathcal{M} \times \Theta$  under the Lévy metric and  $\Theta$  is compact, the maximum theorem (e.g., Berge (1963)) implies that the minimizer  $\bar{T}(Q)$  is upper semi-continuous at each  $Q \in \mathcal{M}$  under the Lévy metric. Since the Hellinger distance is always larger than the Lévy metric for any pair of probability measures,  $\bar{T}(Q)$  is also upper semi-continuous under the Hellinger distance.

**Proof of (ii).** Pick arbitrary  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ . From the triangle inequality, (7.1)

$$\sup_{\theta \in \Theta} |E_{Q_n}[g_n(x, \theta)] - E_{P_0}[g(x, \theta)]| \leq \sup_{\theta \in \Theta} |E_{Q_n}[g_n(x, \theta)] - E_{P_0}[g_n(x, \theta)]| + \sup_{\theta \in \Theta} |E_{P_0}[g(x, \theta) \mathbb{I}\{x \notin \mathcal{X}_n\}]|.$$

The first term of (7.1) satisfies

$$\begin{aligned}
& \sup_{\theta \in \Theta} |E_{Q_n} [g_n(x, \theta)] - E_{P_0} [g_n(x, \theta)]| \\
& \leq \sup_{\theta \in \Theta} \left| \int g_n(x, \theta) \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\}^2 \right| + 2 \sup_{\theta \in \Theta} \left| \int g_n(x, \theta) dP_0^{1/2} \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\} \right| \\
& \leq m_n \frac{r^2}{n} + 2 \sqrt{E_{P_0} \left[ \sup_{\theta \in \Theta} |g(x, \theta)|^2 \right]} \frac{r}{\sqrt{n}} = O(n^{-1/2}),
\end{aligned}$$

where the first inequality follows from the triangle inequality, the second inequality follows from  $Q_n \in B_H(P_0, r/\sqrt{n})$  and the Cauchy-Schwarz inequality, and the equality follows from Assumption 3.1 (v) and (vii). The second term of (7.1) satisfies

$$\begin{aligned}
& \sup_{\theta \in \Theta} |E_{P_0} [g(x, \theta) \mathbb{I}\{x \notin \mathcal{X}_n\}]| \\
& \leq \left( \int \sup_{\theta \in \Theta} |g(x, \theta)|^\eta dP_0 \right)^{1/\eta} \left( \int \mathbb{I}\{x \notin \mathcal{X}_n\} dP_0 \right)^{(\eta-1)/\eta} \\
(7.2) \quad & \leq \left( E_{P_0} \left[ \sup_{\theta \in \Theta} |g(x, \theta)|^\eta \right] \right)^{1/\eta} \left( m_n^{-\eta} E_{P_0} \left[ \sup_{\theta \in \Theta} |g(x, \theta)|^\eta \right] \right)^{(\eta-1)/\eta} = o(n^{-1/2}),
\end{aligned}$$

where the first inequality follows from the Hölder inequality, and the second inequality follows from the Markov inequality, and the equality follows from Assumption 3.1 (v) and (vii). Combining these results, we obtain the uniform convergence  $\sup_{\theta \in \Theta} |E_{Q_n} [g_n(x, \theta)] - E_{P_0} [g(x, \theta)]| \rightarrow 0$ . Therefore, from the triangle inequality and  $|E_{Q_n} [g_n(x, \bar{T}_{Q_n})]| = O(n^{-1/2})$  (Lemma 7.6 (i)),

$$|E_{P_0} [g(x, \bar{T}_{Q_n})]| \leq |E_{P_0} [g(x, \bar{T}_{Q_n})] - E_{Q_n} [g_n(x, \bar{T}_{Q_n})]| + |E_{Q_n} [g_n(x, \bar{T}_{Q_n})]| \rightarrow 0.$$

The conclusion follows from Assumption 3.1 (iii).

**Lemma 7.2.** *Suppose that Assumption 3.1 holds. Then for each  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ ,*

$$(7.3) \quad \sqrt{n} (\bar{T}_{Q_n} - \theta_0) = -\sqrt{n} \Sigma^{-1} \int \Lambda_n dQ_n + o(1).$$

**Proof.** The proof is based on Rieder (1994, proofs of Theorems 6.3.4 and Theorem 6.4.5). Pick arbitrary  $r > 0$  and  $Q_n \in B_H(P_0, r/\sqrt{n})$ . Observe that

$$\begin{aligned}
& \left\| dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} (\bar{T}_{Q_n} - \theta_0)' \Lambda_n dQ_n^{1/2} \right\|^2 \\
&= \left\| dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} \psi'_{n, Q_n} \Lambda_n dQ_n^{1/2} \right\|^2 + \left\| \frac{1}{2} (\bar{T}_{Q_n} - \theta_0 - \psi_{n, Q_n})' \Lambda_n dQ_n^{1/2} \right\|^2 \\
&\quad + \left\{ \int \left( dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} \psi'_{n, Q_n} \Lambda_n dQ_n^{1/2} \right) \Lambda'_n dQ_n^{1/2} \right\} (\bar{T}_{Q_n} - \theta_0 - \psi_{n, Q_n}) \\
(7.4) \quad &= \left\| dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} \psi'_{n, Q_n} \Lambda_n dQ_n^{1/2} \right\|^2 + \left\| \frac{1}{2} (\bar{T}_{Q_n} - \theta_0 - \psi_{n, Q_n})' \Lambda_n dQ_n^{1/2} \right\|^2,
\end{aligned}$$

where the second equality follows from

$$\begin{aligned}
& \int \left\{ dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} \psi'_{n, Q_n} \Lambda_n dQ_n^{1/2} \right\} \Lambda'_n dQ_n^{1/2} \\
&= \int \Lambda'_n \left\{ dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} \right\} dQ_n^{1/2} + \frac{1}{2} \psi'_{n, Q_n} \int \Lambda_n \Lambda'_n dQ_n = 0.
\end{aligned}$$

The left hand side of (7.4) satisfies

$$\begin{aligned}
& \left\| dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} (\bar{T}_{Q_n} - \theta_0)' \Lambda_n dQ_n^{1/2} \right\| \\
&\leq \left\| dQ_n^{1/2} - d\bar{P}_{\bar{T}_{Q_n}, Q_n}^{1/2} \right\| + o(|\bar{T}_{Q_n} - \theta_0|) + o(n^{-1/2}) \\
&\leq \left\| dQ_n^{1/2} - d\bar{P}_{\theta_0 + \psi_{n, Q_n}, Q_n}^{1/2} \right\| + o(|\bar{T}_{Q_n} - \theta_0|) + o(n^{-1/2}) \\
(7.5) \quad &\leq \left\| dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} \psi'_{n, Q_n} \Lambda_n dQ_n^{1/2} \right\| + o(|\bar{T}_{Q_n} - \theta_0|) + o(|\psi_{n, Q_n}|) + o(n^{-1/2}),
\end{aligned}$$

where the first inequality follows from the triangle inequality and Lemma 7.3 (i), the second inequality follows from  $\bar{T}_{Q_n} = \arg \min_{\theta \in \Theta} \left\| dQ_n^{1/2} - d\bar{P}_{\theta, Q_n}^{1/2} \right\|$ , and the third inequality follows from the triangle inequality and Lemma 7.3 (ii). From (7.4) and (7.5),

$$\begin{aligned}
& \left\| \left\| dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} \psi'_{n, Q_n} \Lambda_n dQ_n^{1/2} \right\|^2 + \left\| \frac{1}{2} (\bar{T}_{Q_n} - \theta_0 - \psi_{n, Q_n})' \Lambda_n dQ_n^{1/2} \right\|^2 \right\|^{1/2} \\
&\leq \left\| dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} \psi'_{n, Q_n} \Lambda_n dQ_n^{1/2} \right\| + o(|\bar{T}_{Q_n} - \theta_0|) + o(|\psi_{n, Q_n}|) + o(n^{-1/2}).
\end{aligned}$$

This implies

$$\begin{aligned}
& o(|\bar{T}_{Q_n} - \theta_0|) + o(|\psi_{n, Q_n}|) + o(n^{-1/2}) \\
(7.6) \quad &\geq \sqrt{\frac{1}{4} (\bar{T}_{Q_n} - \theta_0 - \psi_{n, Q_n})' \int \Lambda_n \Lambda'_n dQ_n (\bar{T}_{Q_n} - \theta_0 - \psi_{n, Q_n})} \geq C |\bar{T}_{Q_n} - \theta_0 - \psi_{n, Q_n}|,
\end{aligned}$$

for all  $n$  large enough, where the second inequality follows from Lemma 7.5 (i) and Assumption 3.1 (vi).

We now analyze  $\psi_{n,Q_n}$ . From the definition of  $\psi_{n,Q_n}$ ,

$$(7.7) \quad \begin{aligned} \psi_{n,Q_n} &= -2 \left\{ \left( \int \Lambda_n \Lambda_n' dQ_n \right)^{-1} - \Sigma^{-1} \right\} \int \Lambda_n \left\{ dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} \right\} dQ_n^{1/2} \\ &\quad - 2\Sigma^{-1} \int \Lambda_n \left\{ dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} \right\} dQ_n^{1/2}. \end{aligned}$$

From this and Lemma 7.5 (i), the first term of (7.7) is  $o(n^{-1/2})$ . The second term of (7.7) satisfies

$$\begin{aligned} &-2\Sigma^{-1} \int \Lambda_n \left\{ dQ_n^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} \right\} dQ_n^{1/2} \\ &= -2\Sigma^{-1} G' \Omega^{-1} \left( \int g_n(x, \theta_0) g_n(x, \theta_0)' dQ_n \right) \gamma_n(\theta_0, Q_n) \\ &\quad + 2\Sigma^{-1} G' \Omega^{-1} \left( \int \frac{\gamma_n(\theta_0, Q_n)' g_n(x, \theta_0)}{1 + \gamma_n(\theta_0, Q_n)' g_n(x, \theta_0)} g_n(x, \theta_0) g_n(x, \theta_0)' dQ_n \right) \gamma_n(\theta_0, Q_n) \\ &= -\Sigma^{-1} G' \Omega^{-1} \left\{ \int g_n(x, \theta_0) dQ_n + \frac{1}{2} \int \varrho_n(x, \theta_0, Q_n) g_n(x, \theta_0) dQ_n \right\} + o(n^{-1/2}) \\ &= -\Sigma^{-1} \int \Lambda_n dQ_n + o(n^{-1/2}), \end{aligned}$$

where the first equality follows from (7.8), the second equality follows from (7.9) and Lemma 7.5, and the third equality follows from Lemma 7.5. Therefore,

$$\sqrt{n} \psi_{n,Q_n} = -\sqrt{n} \Sigma^{-1} \int \Lambda_n dQ_n + o(1),$$

which also implies  $|\psi_{n,Q_n}| = O(n^{-1/2})$  (by Lemma 7.5 (i)). Combining this with (7.6),

$$\sqrt{n} (\bar{T}_{Q_n} - \theta_0) = \sqrt{n} \psi_{n,Q_n} + o(\sqrt{n} |\bar{T}_{Q_n} - \theta_0|) + o(1).$$

By solving this equation for  $\sqrt{n} (\bar{T}_{Q_n} - \theta_0)$ , the conclusion is obtained.

**Lemma 7.3.** *Suppose that Assumption 3.1 holds. Then for each  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ ,*

$$\begin{aligned} \text{(i):} & \left\| d\bar{P}_{\bar{T}_{Q_n}, Q_n}^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} (\bar{T}_{Q_n} - \theta_0)' \Lambda_n dQ_n^{1/2} \right\| = o(|\bar{T}_{Q_n} - \theta_0|) + o(n^{-1/2}), \\ \text{(ii):} & \left\| d\bar{P}_{\theta_0 + \psi_{n,Q_n}, Q_n}^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} \psi_{n,Q_n}' \Lambda_n dQ_n^{1/2} \right\| = o(|\psi_{n,Q_n}|) + o(n^{-1/2}). \end{aligned}$$

**Proof of (i).** From the convex duality of partially finite programming (Borwein and Lewis (1993)), the Radon-Nikodym derivative  $d\bar{P}_{\theta, Q}/dQ$  is written as

$$(7.8) \quad \frac{d\bar{P}_{\theta, Q}}{dQ} = \frac{1}{(1 + \gamma_n(\theta, Q)' g_n(x, \theta))^2},$$

for each  $n \in \mathbb{N}$ ,  $\theta \in \Theta$ , and  $Q \in \mathcal{M}$ , where  $\gamma_n(\theta, Q)$  solves

$$(7.9) \quad 0 = \int \frac{g_n(x, \theta)}{(1 + \gamma_n(\theta, Q)' g_n(x, \theta))^2} dQ = E_Q [g_n(x, \theta) \{1 - 2\gamma_n(\theta, Q)' g_n(x, \theta) + \varrho_n(x, \theta, Q)\}],$$

with

$$\varrho_n(x, \theta, Q) = \frac{3(\gamma_n(\theta, Q)' g_n(x, \theta))^2 + 2(\gamma_n(\theta, Q)' g_n(x, \theta))^3}{(1 + \gamma_n(\theta, Q)' g_n(x, \theta))^2}.$$

Denote  $t_n = \bar{T}_{Q_n} - \theta_0$ . Pick arbitrary  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ . From the triangle inequality and (7.8),

$$\begin{aligned} & \left\| d\bar{P}_{\bar{T}_{Q_n}, Q_n}^{1/2} - d\bar{P}_{\theta_0, Q_n}^{1/2} + \frac{1}{2} t_n' \Lambda_n dQ_n^{1/2} \right\| \\ & \leq \left\| \left\{ \gamma_n(\theta_0, Q_n)' g_n(x, \theta_0) - \gamma_n(\bar{T}_{Q_n}, Q_n)' g_n(x, \bar{T}_{Q_n}) \right\} dQ_n^{1/2} + \frac{1}{2} t_n' \Lambda_n dQ_n^{1/2} \right\| \\ & \quad + \left\| \left\{ \gamma_n(\theta_0, Q_n)' g_n(x, \theta_0) - \gamma_n(\bar{T}_{Q_n}, Q_n)' g_n(x, \bar{T}_{Q_n}) \right\} \right. \\ & \quad \left. \times \left\{ \frac{1}{(1 + \gamma_n(\bar{T}_{Q_n}, Q_n)' g_n(x, \bar{T}_{Q_n})) (1 + \gamma_n(\theta_0, Q_n)' g_n(x, \theta_0))} - 1 \right\} dQ_n^{1/2} \right\| = T_1 + T_2. \end{aligned}$$

For  $T_2$ , Lemmas 7.5 and 7.6 imply  $T_2 = o(n^{-1/2})$ . For  $T_1$ , the triangle inequality and (7.9) yield

$$\begin{aligned} T_1 & \leq \left\| \left\{ \begin{aligned} & -\frac{1}{2} E_{Q_n} [g_n(x, \bar{T}_{Q_n})]' E_{Q_n} [g_n(x, \bar{T}_{Q_n}) g_n(x, \bar{T}_{Q_n})']^{-1} g_n(x, \bar{T}_{Q_n}) \\ & + \frac{1}{2} E_{Q_n} [g_n(x, \theta_0)]' E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)']^{-1} g_n(x, \theta_0) + \frac{1}{2} t_n' \Lambda_n \end{aligned} \right\} dQ_n^{1/2} \right\| \\ & \quad + \left\| E_{Q_n} [\varrho_n(x, \theta_0, Q_n) g_n(x, \theta_0)]' E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)']^{-1} g_n(x, \theta_0) dQ_n^{1/2} \right\| \\ & \quad + \left\| E_{Q_n} [\varrho_n(x, \bar{T}_{Q_n}, Q_n) g_n(x, \bar{T}_{Q_n})]' E_{Q_n} [g_n(x, \bar{T}_{Q_n}) g_n(x, \bar{T}_{Q_n})']^{-1} g_n(x, \theta_0) dQ_n^{1/2} \right\| \\ & = T_{11} + T_{12} + T_{13}. \end{aligned}$$

Lemmas 7.5 and 7.6 imply that  $T_{12} = o(n^{-1/2})$  and  $T_{13} = o(n^{-1/2})$ . For  $T_{11}$ , expansions of  $g_n(x, \bar{T}_{Q_n})$  around  $\bar{T}_{Q_n} = \theta_0$  yield

$$\begin{aligned}
T_{11} &\leq \left\| -\frac{1}{2} E_{Q_n} [g_n(x, \bar{T}_{Q_n})]' \begin{pmatrix} E_{Q_n} [g_n(x, \bar{T}_{Q_n}) g_n(x, \bar{T}_{Q_n})']^{-1} \\ -E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)']^{-1} \end{pmatrix} g_n(x, \bar{T}_{Q_n}) dQ_n^{1/2} \right\| \\
&\quad + \left\| -\frac{1}{2} E_{Q_n} [g_n(x, \bar{T}_{Q_n})]' E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)']^{-1} \{g_n(x, \bar{T}_{Q_n}) - g_n(x, \theta_0)\} dQ_n^{1/2} \right\| \\
&\quad + \left\| -\frac{1}{2} t'_n \left( \int \frac{\partial g_n(x, \hat{\theta})}{\partial \theta'} dQ_n - G \right)' E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)']^{-1} g_n(x, \theta_0) dQ_n^{1/2} \right\| \\
&\quad + \left\| \frac{1}{2} t'_n G' \left( \Omega^{-1} - E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)']^{-1} \right) g_n(x, \theta_0) dQ_n^{1/2} \right\| \\
&= o(n^{-1/2}) + o(t_n),
\end{aligned}$$

where  $\hat{\theta}$  is a point on the line joining  $\theta_0$  and  $\bar{T}_{Q_n}$ , and the equality follows from Lemmas 7.5 (i) and 7.6 (i).

**Proof of (ii).** Similar to the proof of Part (i) of this lemma.

**Lemma 7.4.** *Suppose that Assumption 3.1 hold. Then for each  $t \in \mathbb{R}^p$ ,*

- (i):  $|E_{P_0} [g_n(x, \theta_0)]| = o(n^{-1/2})$ ,  $|E_{P_0} [g_n(x, \theta_n)]| = O(n^{-1/2})$ ,  $|E_{P_0} [g_n(x, \theta_n) g_n(x, \theta_n)'] - \Omega| = o(1)$ , and  $|E_{P_0} [\partial g_n(x, \theta_n) / \partial \theta'] - G| = o(1)$ ,
- (ii):  $\gamma_n(\theta_n, P_0) = \arg \max_{\gamma \in \mathbb{R}^m} - \int \frac{1}{(1+\gamma' g_n(x, \theta_n))} dP_0$  exists for all  $n$  large enough,  $|\gamma_n(\theta_n, P_0)| = O(n^{-1/2})$ , and  $\sup_{x \in \mathcal{X}} |\gamma_n(\theta_n, P_0)' g_n(x, \theta_n)| = o(1)$ .

**Proof of (i). Proof of the first statement.** The same argument as (7.2) with Assumption 3.1 (iii) yields the conclusion.

**Proof of the second statement.** Pick an arbitrary  $t \in \mathbb{R}^p$ . From the triangle inequality,

$$(7.10) \quad |E_{P_0} [g_n(x, \theta_n)]| \leq |E_{P_0} [g(x, \theta_n) \mathbb{I}\{x \notin \mathcal{X}_n\}]| + |E_{P_0} [g(x, \theta_n)]|.$$

By the same argument as (7.2) and  $E_{P_0} [|g(x, \theta_n)|^\eta] < \infty$  (from Assumption 3.1 (v)), the first term of (7.10) is  $o(n^{-1/2})$ . The second term of (7.10) satisfies

$$|E_{P_0} [g(x, \theta_n)]| \leq E_{P_0} \left[ \sup_{\theta \in \mathcal{N}} \left| \frac{\partial g(x, \theta)}{\partial \theta'} \right| \right] \left| \frac{t}{\sqrt{n}} \right| = O(n^{-1/2}),$$

for all  $n$  large enough, where the inequality follows from a Taylor expansion around  $t = 0$  and Assumption 3.1 (iii), and the equality follows from Assumption 3.1 (v). Combining these results, the conclusion is obtained.



**Proof of the third statement.** Pick an arbitrary  $t \in \mathbb{R}^p$ . From the triangle inequality,

$$\begin{aligned} & |E_{P_0} [g_n(x, \theta_n) g_n(x, \theta_n)'] - \Omega| \\ \leq & |E_{P_0} [g_n(x, \theta_n) g_n(x, \theta_n)'] - E_{P_0} [g(x, \theta_n) g(x, \theta_n)']| + |E_{P_0} [g(x, \theta_n) g(x, \theta_n)'] - \Omega|. \end{aligned}$$

The first term is  $o(n^{-1/2})$  by the same argument as (7.2) and the second term converges to zero by the continuity of  $g(x, \theta)$  at  $\theta_0$ .

**Proof of the fourth statement.** Similar to the proof of the third statement.

**Proof of (ii).** Pick an arbitrary  $t \in \mathbb{R}^p$ . Let  $\Gamma_n = \{\gamma \in \mathbb{R}^m : |\gamma| \leq a_n\}$  with a positive sequence  $\{a_n\}_{n \in \mathbb{N}}$  satisfying  $a_n m_n \rightarrow 0$  and  $a_n n^{1/2} \rightarrow \infty$ . Observe that

$$(7.11) \quad \sup_{\gamma \in \Gamma_n, x \in \mathcal{X}, \theta \in \Theta} |\dot{\gamma}' g_n(x, \theta)| \leq a_n m_n \rightarrow 0.$$

Since  $R_n(P_0, \theta_n, \gamma)$  is twice continuously differentiable with respect to  $\gamma$  and  $\Gamma_n$  is compact,  $\tilde{\gamma} = \arg \max_{\gamma \in \Gamma_n} R_n(P_0, \theta_n, \gamma)$  exists for each  $n \in \mathbb{N}$ . A Taylor expansion around  $\tilde{\gamma} = 0$  yields

$$\begin{aligned} -1 &= R_n(P_0, \theta_n, 0) \leq R_n(P_0, \theta_n, \tilde{\gamma}) = -1 + \tilde{\gamma}' E_{P_0} [g_n(x, \theta_n)] - \tilde{\gamma}' E_{P_0} \left[ \frac{g_n(x, \theta_n) g_n(x, \theta_n)'}{(1 + \dot{\gamma}' g_n(x, \theta_n))^3} \right] \tilde{\gamma} \\ &\leq -1 + \tilde{\gamma}' E_{P_0} [g_n(x, \theta_n)] - C \tilde{\gamma}' E_{P_0} [g_n(x, \theta_n) g_n(x, \theta_n)'] \tilde{\gamma} \\ (7.12) \quad &\leq -1 + |\tilde{\gamma}| |E_{P_0} [g_n(x, \theta_n)]| - C |\tilde{\gamma}|^2, \end{aligned}$$

for all  $n$  large enough, where  $\dot{\gamma}$  is a point on the line joining 0 and  $\tilde{\gamma}$ , the second inequality follows from (7.11), and the last inequality follows from Lemma 7.4 (i) and Assumption 3.1 (vi). Thus, Lemma 7.4 (i) implies

$$(7.13) \quad C |\tilde{\gamma}| \leq |E_{P_0} [g_n(x, \theta_n)]| = O(n^{-1/2}).$$

From  $a_n n^{1/2} \rightarrow \infty$ ,  $\tilde{\gamma}$  is an interior point of  $\Gamma_n$  and satisfies the first-order condition  $\partial R_n(Q_n, \theta_0, \tilde{\gamma}) / \partial \gamma = 0$  for all  $n$  large enough. Since  $R_n(Q_n, \theta_0, \gamma)$  is concave in  $\gamma$  for all  $n$  large enough,  $\tilde{\gamma} = \arg \max_{\gamma \in \mathbb{R}^m} R_n(P_0, \theta_n, \gamma)$  for all  $n$  large enough and the first statement is obtained. Thus, the second statement is obtained from (7.13). The third statement follows from (7.13) and Assumption 3.1 (vii).

**Lemma 7.5.** *Suppose that Assumption 3.1 holds. Then for each  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ ,*

- (i):  $|E_{Q_n} [g_n(x, \theta_0)]| = O(n^{-1/2})$ , and  $|E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)'] - \Omega| = o(1)$ ,
- (ii):  $\gamma_n(\theta_0, Q_n) = \arg \max_{\gamma \in \mathbb{R}^m} - \int \frac{1}{(1 + \dot{\gamma}' g_n(x, \theta_0))} dQ_n$  exists for all  $n$  large enough, and  $|\gamma_n(\theta_0, Q_n)| = O(n^{-1/2})$ , and  $\sup_{x \in \mathcal{X}} |\gamma_n(\theta_0, Q_n)' g_n(x, \theta_0)| = o(1)$ .

**Proof of (i). Proof of the first statement.** Pick any  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ . We have

$$\begin{aligned}
& |E_{Q_n} [g_n(x, \theta_0)]| \\
& \leq \left| \int g_n(x, \theta_0) \{dQ_n - dP_0\} \right| + |E_{P_0} [g_n(x, \theta_0)]| \\
& \leq \left| \int g_n(x, \theta_0) \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\}^2 \right| + 2 \left| \int g_n(x, \theta_0) dP_0^{1/2} \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\} \right| + o(n^{-1/2}) \\
& \leq m_n \frac{r^2}{n} + 2E_{P_0} \left[ |g(x, \theta_0)|^2 \right] \frac{r}{\sqrt{n}} + o(n^{-1/2}) = O(n^{-1/2}),
\end{aligned}$$

where the first and second inequalities follow from the triangle inequality and Lemma 7.4 (i), the third inequality follows from the Cauchy-Schwarz inequality and  $Q_n \in B_H(P_0, r/\sqrt{n})$ , and the equality follows from Assumption 3.1 (v) and (vii).

**Proof of the second statement.** Pick arbitrary  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ . From the triangle inequality,

$$\begin{aligned}
(7.14) \quad & |E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)'] - \Omega| \\
& \leq |E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)'] - E_{P_0} [g_n(x, \theta_0) g_n(x, \theta_0)']| + |E_{P_0} [g(x, \theta_0) g(x, \theta_0)' \mathbb{I}\{x \notin \mathcal{X}_n\}]|.
\end{aligned}$$

The first term of the RHS of (7.14) satisfies

$$\begin{aligned}
& |E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)'] - E_{P_0} [g_n(x, \theta_0) g_n(x, \theta_0)']| \\
& \leq \left| \int g_n(x, \theta_0) g_n(x, \theta_0)' \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\}^2 \right| + 2 \left| \int g_n(x, \theta_0) g_n(x, \theta_0)' dP_0^{1/2} \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\} \right| \\
& \leq m_n^2 \frac{r^2}{n} + 2E_{P_0} \left[ |g(x, \theta_0)|^4 \right] \frac{r}{\sqrt{n}} = o(1),
\end{aligned}$$

where the first inequality follows from the triangle inequality, the second inequality follows from the Cauchy-Schwarz inequality and  $Q_n \in B_H(P_0, r/\sqrt{n})$ , and the equality follows from Assumption 3.1 (v) and (vii). The second term of the RHS of (7.14) satisfies

$$\begin{aligned}
& |E_{P_0} [g(x, \theta_0) g(x, \theta_0)' \mathbb{I}\{x \notin \mathcal{X}_n\}]| \\
& \leq \left( \int |g(x, \theta_0) g(x, \theta_0)'|^{1+\delta} dP_0 \right)^{\frac{1}{1+\delta}} \left( \int \mathbb{I}\{x \notin \mathcal{X}_n\} dP_0 \right)^{\frac{\delta}{1+\delta}} \\
& \leq \left( E_{P_0} \left[ |g(x, \theta_0)|^{2+\delta} \right] \right)^{\frac{1}{1+\delta}} \left( m_n^{-\eta} E_{P_0} [|g(x, \theta_0)|^\eta] \right)^{\frac{\delta}{1+\delta}} = o(1),
\end{aligned}$$

for sufficiently small  $\delta > 0$ , where the first inequality follows from the Hölder inequality, the second inequality follows from the Markov inequality, and the equality follows from Assumption 3.1 (vii).

**Proof of (ii).** Similar to the proof of Lemma 7.4 (ii). Repeat the same argument with  $R_n(Q_n, \theta_0, \gamma)$  instead of  $R_n(P_0, \theta_n, \gamma)$ .

**Lemma 7.6.** *Suppose that Assumption 3.1 holds. Then for each  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ ,*

- (i):  $|E_{Q_n}[g_n(x, \bar{T}_{Q_n})]| = O(n^{-1/2})$ ,  $|E_{Q_n}[g_n(x, \bar{T}_{Q_n})g_n(x, \bar{T}_{Q_n})'] - \Omega| = o(1)$ , and  $|E_{Q_n}[\partial g_n(x, \bar{T}_{Q_n})/\partial \theta'] - G| = o(1)$ ,
- (ii):  $\gamma_n(\bar{T}_{Q_n}, Q_n) = \arg \max_{\gamma \in \mathbb{R}^m} - \int \frac{1}{(1 + \gamma' g_n(x, \bar{T}_{Q_n}))} dQ_n$  exists for all  $n$  large enough,  $|\gamma_n(\bar{T}_{Q_n}, Q_n)| = O(n^{-1/2})$ , and  $\sup_{x \in \mathcal{X}} |\gamma_n(\bar{T}_{Q_n}, Q_n)' g_n(x, \bar{T}_{Q_n})| = o(1)$ .

**Proof of (i). Proof of the first statement.** Pick any  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ .

Define  $\tilde{\gamma} = \frac{E_{Q_n}[g_n(x, \bar{T}_{Q_n})]}{\sqrt{n}|E_{Q_n}[g_n(x, \bar{T}_{Q_n})]|}$ . Since  $|\tilde{\gamma}| = n^{-1/2}$ ,

$$(7.15) \quad \sup_{x \in \mathcal{X}, \theta \in \Theta} |\tilde{\gamma}' g_n(x, \theta)| \leq n^{-1/2} m_n \rightarrow 0.$$

Observe that

$$(7.16) \quad \begin{aligned} & \left| E_{Q_n} \left[ g_n(x, \bar{T}_{Q_n}) g_n(x, \bar{T}_{Q_n})' \right] \right| \\ & \leq \int \sup_{\theta \in \Theta} |g_n(x, \theta)|^2 \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\}^2 + 2 \int \sup_{\theta \in \Theta} |g_n(x, \theta)|^2 dP_0^{1/2} \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\} + E_{P_0} \left[ \sup_{\theta \in \Theta} |g_n(x, \theta)|^2 \right] \\ & \leq m_n^2 \frac{r^2}{n} + 2m_n \sqrt{E_{P_0} \left[ \sup_{\theta \in \Theta} |g(x, \theta)|^2 \right]} \frac{r}{\sqrt{n}} + E_{P_0} \left[ \sup_{\theta \in \Theta} |g(x, \theta)|^2 \right] \leq CE_{P_0} \left[ \sup_{\theta \in \Theta} |g(x, \theta)|^2 \right], \end{aligned}$$

for all  $n$  large enough, where the first inequality follows from the triangle inequality, the second inequality follows from the Cauchy-Schwarz inequality and  $Q_n \in B_H(P_0, r/\sqrt{n})$ , and the last inequality follows from Assumption 3.1 (v) and (vii). Thus, an expansion around  $\tilde{\gamma} = 0$  yields

$$(7.17) \quad \begin{aligned} R_n(Q_n, \bar{T}_{Q_n}, \tilde{\gamma}) &= -1 + \tilde{\gamma}' E_{Q_n}[g_n(x, \bar{T}_{Q_n})] - \tilde{\gamma}' E_{Q_n} \left[ \frac{g_n(x, \bar{T}_{Q_n}) g_n(x, \bar{T}_{Q_n})'}{(1 + \tilde{\gamma}' g_n(x, \bar{T}_{Q_n}))^3} \right] \tilde{\gamma} \\ &\geq -1 + n^{-1/2} |E_{Q_n}[g_n(x, \bar{T}_{Q_n})]| - C\tilde{\gamma}' E_{Q_n} \left[ g_n(x, \bar{T}_{Q_n}) g_n(x, \bar{T}_{Q_n})' \right] \tilde{\gamma} \\ &\geq -1 + n^{-1/2} |E_{Q_n}[g_n(x, \bar{T}_{Q_n})]| - Cn^{-1}, \end{aligned}$$

for all  $n$  large enough, where  $\dot{\gamma}$  is a point on the line joining 0 and  $\tilde{\gamma}$ , the first inequality follows from (7.15), and the second inequality follows from  $\tilde{\gamma}' \tilde{\gamma} = n^{-1}$  and (7.16). From the duality of partially finite programming (Borwein and Lewis (1993)),  $\gamma_n(\bar{T}_{Q_n}, Q_n)$  and  $\bar{T}_{Q_n}$  are written as  $\gamma_n(\bar{T}_{Q_n}, Q_n) =$

$\arg \max_{\gamma \in \mathbb{R}^m} R_n(Q_n, \bar{T}_{Q_n}, \gamma)$  and  $\bar{T}_{Q_n} = \arg \min_{\theta \in \Theta} R_n(Q_n, \theta, \gamma_n(\theta, Q_n))$ . Therefore, from (7.17),

$$(7.18) \quad \begin{aligned} & -1 + n^{-1/2} |E_{Q_n} [g_n(x, \bar{T}_{Q_n})]| - Cn^{-1} \\ & \leq R_n(Q_n, \bar{T}_{Q_n}, \tilde{\gamma}) \leq R_n(Q_n, \bar{T}_{Q_n}, \gamma_n(\bar{T}_{Q_n}, Q_n)) \leq R_n(Q_n, \theta_0, \gamma_n(\theta_0, Q_n)). \end{aligned}$$

By a similar argument to (7.12) combined with  $|\gamma_n(\theta_0, Q_n)| = O(n^{-1/2})$  and  $|E_{Q_n} [g_n(x, \theta_0)]| = O(n^{-1/2})$  (by Lemma 7.5), we have

$$(7.19) \quad R_n(Q_n, \theta_0, \gamma_n(\theta_0, Q_n)) \leq -1 + |\gamma_n(\theta_0, Q_n)| |E_{Q_n} [g_n(x, \theta_0)]| - C |\gamma_n(\theta_0, Q_n)|^2 = -1 + O(n^{-1}).$$

From (7.18) and (7.19), the conclusion follows.

**Proof of the second statement.** Similar to the proof of the second statement of Lemma 7.5 (i).

**Proof of the third statement.** Pick arbitrary  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ . From the triangle inequality,

$$(7.20) \quad \begin{aligned} |E_{Q_n} [\partial g_n(x, \bar{T}_{Q_n}) / \partial \theta'] - G| & \leq |E_{Q_n} [\partial g_n(x, \bar{T}_{Q_n}) / \partial \theta'] - E_{P_0} [\partial g_n(x, \bar{T}_{Q_n}) / \partial \theta']| \\ & \quad + |E_{P_0} [\mathbb{I}\{x \notin \mathcal{X}_n\} \partial g(x, \bar{T}_{Q_n}) / \partial \theta']| + |E_{P_0} [\partial g(x, \bar{T}_{Q_n}) / \partial \theta'] - G|. \end{aligned}$$

The first term of (7.20) satisfies

$$\begin{aligned} & |E_{Q_n} [\partial g_n(x, \bar{T}_{Q_n}) / \partial \theta'] - E_{P_0} [\partial g_n(x, \bar{T}_{Q_n}) / \partial \theta']| \\ & \leq \left| \int \partial g_n(x, \bar{T}_{Q_n}) / \partial \theta' \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\}^2 \right| + 2 \left| \int \partial g_n(x, \bar{T}_{Q_n}) / \partial \theta' dP_0^{1/2} \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\} \right| \\ & \leq \sup_{x \in \mathcal{X}_n, \theta \in \mathcal{N}} |\partial g_n(x, \theta) / \partial \theta'| \frac{r^2}{n} + 2E_{P_0} \left[ \sup_{\theta \in \mathcal{N}} |\partial g_n(x, \theta) / \partial \theta'|^2 \right] \frac{r}{\sqrt{n}} = o(1), \end{aligned}$$

where the first inequality follows from the triangle inequality, the second inequality follows from the Cauchy-Schwarz inequality, and the equality follows from Assumption 3.1 (v) and (vii). The second term of (7.20) is  $o(1)$  by the same argument as (7.2). The third term of (7.20) is  $o(1)$  by the continuity of  $\partial g(x, \theta) / \partial \theta'$  at  $\theta_0$  and Lemma 7.1 (ii). Therefore, the conclusion is obtained.

**Proof of (ii).** Similar to the proof of Lemma 7.4 (ii). Repeat the same argument with  $R_n(Q_n, \bar{T}_{Q_n}, \gamma)$  instead of  $R_n(P_0, \theta_n, \gamma)$ .

**Lemma 7.7.** *Suppose that Assumption 3.1 holds. Then for each sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$  and  $r > 0$ ,  $\bar{T}_{P_n} \xrightarrow{P} \theta_0$  under  $Q_n$ .*

**Proof.** Similar to the proof of Lemma 7.1 (i).

**Lemma 7.8.** *Suppose that Assumption 3.1 holds . Then for each  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ ,*

$$\begin{aligned}\sqrt{n}(\bar{T}_{P_n} - \theta_0) &= -\sqrt{n}\Sigma^{-1} \int \Lambda_n dP_n + o_p(1) \quad \text{under } Q_n, \\ \sqrt{n}(\bar{T}_{P_n} - \bar{T}_{Q_n}) &\xrightarrow{d} N(0, \Sigma^{-1}) \quad \text{under } Q_n.\end{aligned}$$

**Proof.** The proof of the first statement is similar to that of Lemma 7.2 (replace  $Q_n$  with  $P_n$  and use Lemmas 7.9 and 7.10 instead of Lemmas 7.5 and 7.6). For the second statement, Lemma 7.2 and the first statement imply

$$\sqrt{n}(\bar{T}_{P_n} - \bar{T}_{Q_n}) = -\Sigma^{-1}G'\Omega^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \{g_n(x_i, \theta_0) - E_{Q_n}[g_n(x, \theta_0)]\} + o_p(1),$$

under  $Q_n$ . Thus, it is sufficient to check that we can apply a central limit theorem to the triangular array  $\{g_n(x_i, \theta_0)\}_{1 \leq i \leq n, n}$ . Observe that

$$\begin{aligned}& E_{Q_n} \left[ |g_n(x, \theta_0)|^{2+\epsilon} \right] \\ &= \int |g_n(x, \theta_0)|^{2+\epsilon} \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\}^2 + 2 \int |g_n(x, \theta_0)|^{2+\epsilon} dP_0^{1/2} \left\{ dQ_n^{1/2} - dP_0^{1/2} \right\} + E_{P_0} \left[ |g_n(x, \theta_0)|^{2+\epsilon} \right] \\ &\leq m_n^{2+\epsilon} \frac{r^2}{n} + 2m_n^{1+\epsilon} E_{P_0} \left[ |g(x, \theta_0)|^2 \right] \frac{r}{\sqrt{n}} + E_{P_0} \left[ |g(x, \theta_0)|^{2+\epsilon} \right] < \infty,\end{aligned}$$

for all  $n$  large enough, where the first inequality follows from the Cauchy-Schwarz inequality, and the second inequality follows from Assumption 3.1 (v) and (vii). Therefore, the conclusion is obtained.

**Lemma 7.9.** *Suppose that Assumption 3.1 holds. Then for each  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ , the followings hold under  $Q_n$ :*

- (i):  $|E_{P_n}[g_n(x, \theta_0)]| = O_p(n^{-1/2})$ ,  $|E_{P_n}[g_n(x, \theta_0)g_n(x, \theta_0)'] - \Omega| = o_p(1)$ ,
- (ii):  $\gamma_n(\theta_0, P_n) = \arg \max_{\gamma \in \mathbb{R}^m} - \int \frac{1}{(1+\gamma'g_n(x, \theta_0))} dP_n$  exists a.s. for all  $n$  large enough,  $|\gamma_n(\theta_0, P_n)| = O_p(n^{-1/2})$ , and  $\sup_{x \in \mathcal{X}} |\gamma_n(\theta_0, P_n)'g_n(x, \theta_0)| = o_p(1)$ .

**Proof of (i). Proof of the first statement.** From the triangle inequality,

$$|E_{P_n}[g_n(x, \theta_0)]| \leq |E_{P_n}[g_n(x, \theta_0)] - E_{Q_n}[g_n(x, \theta_0)]| + |E_{Q_n}[g_n(x, \theta_0)]|.$$

The first term is  $O_p(n^{-1/2})$  by the central limit theorem for the triangular array  $\{g_n(x_i, \theta_0)\}_{1 \leq i \leq n, n}$ .

The second term is  $O(n^{-1/2})$  by Lemma 7.5 (i).

**Proof of the second statement.** From the triangle inequality,

$$\begin{aligned} & \left| E_{P_n} [g_n(x, \theta_0) g_n(x, \theta_0)' - \Omega] \right| \\ & \leq \left| E_{P_n} [g_n(x, \theta_0) g_n(x, \theta_0)'] - E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)'] \right| + \left| E_{Q_n} [g_n(x, \theta_0) g_n(x, \theta_0)'] - \Omega \right|. \end{aligned}$$

From a law of large numbers, the first term is  $o_p(1)$ . From Lemma 7.5 (i), the second term is  $o(1)$ .

**Proof of (ii).** Similar to the proof of Lemma 7.4 (ii) except using Lemma 7.9 (i) instead of Lemma 7.4 (i).

**Lemma 7.10.** *Suppose that Assumption 3.1 holds. Then for each  $r > 0$  and sequence  $Q_n \in B_H(P_0, r/\sqrt{n})$ , the followings hold under  $Q_n$ :*

- (i):  $|E_{P_n} [g_n(x, \bar{T}_{P_n})]| = O_p(n^{-1/2})$ ,  $|E_{P_n} [g_n(x, \bar{T}_{P_n}) g_n(x, \bar{T}_{P_n})'] - \Omega| = O_p(n^{-1/2})$ , and  $|E_{P_n} [\partial g_n(x, \bar{T}_{P_n}) / \partial \theta'] - G| = o_p(1)$ ,
- (ii):  $\gamma_n(\bar{T}_{P_n}, P_n) = \arg \max_{\gamma \in \mathbb{R}^m} - \int \frac{1}{(1 + \gamma' g_n(x, \bar{T}_{P_n}))} dP_n$  exists a.s. for all  $n$  large enough,  $|\gamma_n(\bar{T}_{P_n}, P_n)| = O_p(n^{-1/2})$ , and  $\sup_{x \in \mathcal{X}} |\gamma_n(\bar{T}_{P_n}, P_n)' g_n(x, \bar{T}_{P_n})| = o_p(1)$ .

**Proof of (i).** Similar to the proof of Lemma 7.6 (i).

**Proof of (ii).** Similar to the proof of Lemma 7.6 (ii).

## REFERENCES

- BERAN, R. (1977): “Minimum Hellinger distance estimates for parametric models,” *Annals of Statistics*, 5, 445–463.
- (1984): “Minimum distance procedures,” in *Handbook of Statistics*, ed. by P. Krishnaiah, and P. Sen. Elsevier Science, pp. 741–754.
- BERGE, C. (1963): *Topological Spaces*. Dover.
- BICKEL, P., C. KLASSEN, Y. RITOV, AND J. WELLNER (1993): *Efficient and Adaptive Estimation for Semiparametric Models*. Baltimore: Johns Hopkins Press.
- BICKEL, P. J. (1981): “Quelques aspects de la statistique robuste,” in *Ecole d’Eté de Probabilités de Saint Flour IX 1979*, ed. by P. Hennequin. Springer, pp. 1–72.
- BORWEIN, J. M., AND A. S. LEWIS (1993): “Partially-finite programming in  $L_1$  and the existence of maximum entropy estimates,” *SIAM Journal of Optimization*, 3, 248–267.
- DONOHO, D., AND R. LIU (1988): “The “automatic” robustness of minimum distance functionals,” *Annals of Statistics*, 16, 552–586.
- HALL, P., AND J. L. HOROWITZ (1996): “Bootstrap Critical Values for Tests Based on Generalized Method of Moments Estimators,” *Econometrica*, 64, 891–916.
- HANSEN, L. P. (1982): “Large Sample Properties of Generalized Methods of Moments Estimators,” *Econometrica*, 50, 1029–1054.
- IMBENS, G. W., R. H. SPADY, AND P. JOHNSON (1998): “Information Theoretic Approaches to Inference in Moment Condition Models,” *Econometrica*, 66, 333–357.
- KITAMURA, Y. (1998): “Comparing Misspecified Dynamic Econometric Models Using Nonparametric Likelihood,” Working Paper, Department of Economics, University of Wisconsin.
- (2001): “Asymptotic optimality of empirical likelihood for testing moment restrictions,” *Econometrica*, 69, 1661–1672.
- (2002): “A Likelihood-based Approach to the Analysis of a Class of Nested and Non-nested Models,” Working Paper, Department of Economics, University of Pennsylvania.
- (2006): “Empirical Likelihood Methods in Econometrics: Theory and Practice,” in *Advances in Economics and Econometrics: Theory and Applications, Ninth World Congress*, ed. by R. Blundell, W. K. Newey, and T. Persson. Cambridge: Cambridge University Press, forthcoming.
- KITAMURA, Y., AND M. STUTZER (1997): “An Information Theoretic Alternative to Generalized Method of Moments Estimation,” *Econometrica*, 65(4), 861–874.
- NEWAY, W. K. (1990): “Semiparametric Efficiency Bounds,” *Journal of Applied Econometrics*, 5, 99–135.
- NEWAY, W. K., AND R. J. SMITH (2004): “Higher order properties of GMM and Generalized Empirical Likelihood Estimators,” *Econometrica*, 72, 219–255.
- POLLARD, D. (2002): *A User’s Guide to Measure Theoretic Probability*. Cambridge.
- REISS, R.-D. (1989): *Approximate Distributions of Order Statistics*. Springer-Verlag.
- RIEDER, H. (1994): *Robust Asymptotic Statistics*. Springer-Verlag.
- ROCKAFELLER, R. (1970): *Convex Analysis*. Princeton University Press.

- SCHENNACH, S. M. (2007): "Point estimation with exponentially tilted empirical likelihood," *Annals of Statistics*, 35, 634–672.
- SMITH, R. J. (1997): "Alternative semi-parametric likelihood approaches to generalized method of moments estimation," *Economic Journal*, 107, 503–519.
- VAN DER VAART, A. (1998): *Asymptotic Statistics*. Cambridge University Press.
- WHITE, H. (1982): "Maximum likelihood estimation of misspecified models," *Econometrica*, 50, 1–25.
- ZHANG, T. (2006): "From  $\varepsilon$ -entropy to KL-entropy: Analysis of Minimum Information Complexity Density Estimation," *The Annals of Statistics*, 34, 2180–2210.

COWLES FOUNDATION FOR RESEARCH IN ECONOMICS, YALE UNIVERSITY, NEW HAVEN, CT-06520.

*E-mail address:* `yuichi.kitamura@yale.edu`

COWLES FOUNDATION FOR RESEARCH IN ECONOMICS, YALE UNIVERSITY, NEW HAVEN, CT-06520.

*E-mail address:* `taisuke.otsu@yale.edu`

DEPARTMENT OF ECONOMICS, YALE UNIVERSITY, NEW HAVEN, CT-06520.

*E-mail address:* `kirill.evdokimov@yale.edu`