

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

3  
 4                   **BY YUICHI KITAMURA, TAISUKE OTSU, AND KIRILL EVDOKIMOV<sup>1</sup>**

5                   This paper is concerned with robust estimation under moment restrictions. A mo-  
 6                   ment restriction model is semiparametric and distribution-free; therefore it imposes  
 7                   mild assumptions. Yet it is reasonable to expect that the probability law of observations  
 8                   may have some deviations from the ideal distribution being modeled, due to various  
 9                   factors such as measurement errors. It is then sensible to seek an estimation proce-  
 10                  dure that is robust against slight perturbation in the probability measure that generates  
 11                  observations. This paper considers local deviations within shrinking topological neigh-  
 12                  borhoods to develop its large sample theory, so that both bias and variance matter  
 13                  asymptotically. The main result shows that there exists a computationally convenient  
 14                  estimator that achieves optimal minimax robust properties. It is semiparametrically ef-  
 15                  ficient when the model assumption holds, and, at the same time, it enjoys desirable  
 16                  robust properties when it does not.

16                  **KEYWORDS:** Asymptotic Minimax Theorem, Hellinger distance, semiparametric ef-  
 17                  ficiency.

18  
 19  
 20   **1. INTRODUCTION**

21                  **CONSIDER A PROBABILITY MEASURE**  $P_0 \in \mathcal{M}$ , where  $\mathcal{M}$  is the set of all proba-  
 22                  bility measures on the Borel  $\sigma$ -field  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$  of  $\mathcal{X} \subseteq \mathbb{R}^d$ . Let  $g : \mathcal{X} \times \Theta \rightarrow \mathbb{R}^m$   
 23                  be a vector of functions parameterized by a  $p$ -dimensional vector  $\theta$  which re-  
 24                  sides in  $\Theta \subset \mathbb{R}^p$ . The function  $g$  satisfies

25  
 26                  (1.1)       $E_{P_0}[g(x, \theta_0)] = \int g(x, \theta_0) dP_0 = 0, \quad \theta_0 \in \Theta.$

27  
 28  
 29                  The moment condition model (1.1) is semiparametric and distribution-free;  
 30                  therefore it imposes mild assumptions. Nevertheless, it is reasonable to expect  
 31                  that the probability law of observations may have some deviations from the re-  
 32                  striction under the moment condition model. It is then sensible to seek for es-  
 33                  timation and testing procedures that are robust against slight perturbations in  
 34                  the observed data, or more formally, perturbations in the probability measure  
 35                  that generates observations. This notion of robustness can be illustrated as fol-  
 36                  lows. Let a functional  $\theta(P)$ ,  $P \in \mathcal{M}$  solve the moment condition model (1.1), in  
 37                  the sense that  $\theta_0 = \theta(P_0)$ . Suppose, however, that observations  $x_1, \dots, x_n$  are

38  
 39                  <sup>1</sup>We are grateful to a co-editor and two anonymous referees for helpful remarks and sugges-  
 40                  tions. We thank participants at the CIREQ Conference on GMM, the 2007 Winter Meetings of  
 41                  the Econometric Society, the 2007 Netherlands Econometrics Study Group Annual Conference,  
 42                  and seminars at Boston University, Chicago Booth, Harvard-MIT, LSE, NYU, Ohio State, Seoul  
 43                  National University, Texas A&M, the University of Tokyo, Vanderbilt, and Wisconsin for valuable  
 44                  comments. We acknowledge financial support from the National Science Foundation via Grants  
 SES-0241770, SES-0551271, and SES-0851759 (Kitamura) and SES-0720961 (Otsu).

1 not drawn according to  $P_0$ , but according to its “perturbed” version  $P$  instead. 1  
 2 This can be attributed to various factors, including measurement errors or data 2  
 3 contamination. These are imminent and realistic concerns in applications. The 3  
 4 goal of robust estimation is to obtain an estimator  $\bar{\theta} = \bar{\theta}(x_1, \dots, x_n)$  that is not 4  
 5 sensitive to such perturbations, so that the deviation of the estimated value  $\bar{\theta}$  5  
 6 from the parameter value of interest  $\theta_0 = \theta(P_0)$  remains stable. Decompose 6  
 7 the deviation as 7

$$8 \quad (1.2) \quad \bar{\theta} - \theta_0 = [\bar{\theta} - \theta(P)] + [\theta(P) - \theta(P_0)]. \quad 8$$

10 In the asymptotic mean squared error (MSE) calculation presented below, the 10  
 11 expectation of the square of the term in the first square bracket contributes to 11  
 12 the variance of the estimator, whereas the second corresponds to the bias. An 12  
 13 estimator that achieves small MSE *uniformly* in  $P$  over a neighborhood of  $P_0$  is 13  
 14 desirable. 14

15 Asymptotic theory of robust estimation when the model is parametric has 15  
 16 been considered extensively in the literature; see [Rieder \(1994\)](#) for a compre- 16  
 17 hensive survey. In a pioneering paper, [Beran \(1977\)](#) discussed “robust and effi- 17  
 18 cient” estimation of parametric models. Suppose  $P_\theta$ ,  $\theta \in \Theta \subset \mathbb{R}^k$  is a paramet- 18  
 19 ric family of probability measures. Observations are drawn from a probability 19  
 20 law  $P$ , which may not be a member of the parametric family. Let  $p_\theta$  and  $p$  de- 20  
 21 note the densities associated with the probability measures  $P_\theta$  and  $P$ . It is well 21  
 22 known that the parametric MLE procedure corresponds to minimizing the 22  
 23 objective function  $\rho = \int \log(p/p_\theta)p \, dx$ . Beran pointed out that a small change 23  
 24 in the density  $p$  can lead to a large change in the objective function  $\rho$  (note 24  
 25 the log in  $\rho$ ), implying the nonrobustness of the MLE. He showed that the 25  
 26 parametric minimum Hellinger distance estimator (MHDE) is “robust and ef- 26  
 27 ficient,” in the sense that (i) it has an asymptotic minimax robust property, and 27  
 28 (ii) it is asymptotically efficient when the model assumption is satisfied, that is, 28  
 29 when the sample is generated from  $P_0 = P_{\theta_0}$ , where  $\theta_0$  is the true value of the 29  
 30 parameter of interest. Let  $H(P_\theta, P) = \sqrt{\int (p_\theta^{1/2}(x) - p^{1/2}(x))^2 \, dx}$  denote the 30  
 31 Hellinger distance between  $P_\theta$  and  $P$  (a slightly more general definition of the 31  
 32 Hellinger distance is given in the next section). The MHDE for the parametric 32  
 33 model is 33  
 34

$$35 \quad \hat{\theta} = \operatorname{argmin}_{\theta} H(P_\theta, \hat{P}) \quad 35$$

$$36 \quad = \operatorname{argmin}_{\theta} \int (p_\theta^{1/2}(x) - \hat{p}^{1/2}(x))^2 \, dx, \quad 36$$

37 where  $\hat{p}$  is a nonparametric density estimator, such as a kernel density estima- 37  
 38 tor, for  $P$ , and  $\hat{P}$  is the corresponding estimator for the probability measure 38  
 39 of  $x$ . The MHDE is asymptotically equivalent to MLE and thus efficient if 39  
 40 the model assumption is satisfied. One can replace the Hellinger distance with 40  
 41  
 42  
 43  
 44

1 other divergence measures such as the Kolmogorov–Smirnov distance, which 1  
 2 would make the corresponding minimum divergence estimator even more ro- 2  
 3 bust, but it would incur efficiency loss. The parametric MHDE has been stud- 3  
 4 ied extensively and applied to various models. 4

5 The parametric MHDE has theoretical advantages and excellent finite sam- 5  
 6 ple performance documented by numerous simulation studies, but it has lim- 6  
 7 itations as well. It requires the nonparametric density estimator when at least 7  
 8 some components of  $x$  are continuously distributed. This makes its practi- 8  
 9 cal application inconvenient, and is problematic when  $x$  is high dimensional, 9  
 10 due to the curse of dimensionality. It also necessitates the evaluation of the 10  
 11 integral  $\int (p_{\theta}^{1/2}(x) - \hat{p}^{1/2}(x))^2 dx$ . This involves either numerical integration 11  
 12 or an approximation by an empirical average with inverse density weighting 12  
 13 using a nonparametric density estimator. The former can be hard to imple- 13  
 14 ment, and the latter may have undesirable effects in finite samples. This paper 14  
 15 aims at developing robust methods for moment restriction models, by applying 15  
 16 the MHDE procedure. The resulting estimator is semiparametrically efficient 16  
 17 when the model assumption holds, and, at the same time, it enjoys an optimal 17  
 18 minimax robust property when it does not. The implementation of the estima- 18  
 19 tor is easy. Unlike its parametric predecessor, it requires neither nonparamet- 19  
 20 ric density estimation nor evaluation of integration. 20  
 21

## 22 2. PRELIMINARIES 23

24 The econometrician wishes to estimate the unknown  $\theta_0$  in (1.1). Suppose a 25  
 26 random sample  $\{x_i\}_{i=1}^n$  generated from  $P$  is observed. As discussed in Section 1, 26  
 27 our focus is on robust estimation of  $\theta_0$  when the probability measure  $P$ , from 27  
 28 which the observations are drawn, is a (locally) perturbed version of  $P_0$ , not 28  
 29  $P_0$  itself. There exists an extensive literature concerning the estimation of (1.1) 29  
 30 under the “classical” setting where data are indeed drawn from  $P_0$ . Many es- 30  
 31 timators for  $\theta_0$  are available, including GMM (Hansen (1982)), the empirical 31  
 32 likelihood (EL) estimator, and its variants. This paper is concerned with an 32  
 33 estimator that can be viewed as MHDE applied to the moment restriction model 33  
 34 (1.1). The Hellinger distance between two probability measures is defined as 34  
 35 follows: 35  
 36

37  
 38 DEFINITION 2.1: Let  $P$  and  $Q$  be probability measures with densities  $p$  and 38  
 39  $q$  with respect to a dominating measure  $\nu$ . The Hellinger distance between  $P$  39  
 40 and  $Q$  is then given by 40  
 41

$$42 H(P, Q) = \left\{ \int (p^{1/2} - q^{1/2})^2 d\nu \right\}^{1/2} = \left\{ 2 - 2 \int p^{1/2} q^{1/2} d\nu \right\}^{1/2}. \quad 43$$

44

1 It is often convenient to use the standard notation in the literature that does 1  
2 not explicitly refer to the dominating measure. Then the above definition be- 2  
3 comes 3

$$4 \quad H(P, Q) = \left\{ \int (dP^{1/2} - dQ^{1/2})^2 \right\}^{1/2} = \left\{ 2 - 2 \int dP^{1/2} dQ^{1/2} \right\}^{1/2}. \quad 4$$

5 Here we show some results concerning the Hellinger distance that are useful 5  
6 in understanding the robustness theorems in the next section. 6  
7

8 DEFINITION 2.2: Let  $P$  and  $Q$  be probability measures with densities  $p$  and 8  
9  $q$  with respect to a dominating measure  $\nu$ . The  $\alpha$ -divergence from  $Q$  to  $P$  is 9  
10 given by 10

$$11 \quad I_\alpha(P, Q) = \frac{1}{\alpha(1-\alpha)} \int \left( 1 - \left( \frac{p}{q} \right)^\alpha \right) q \, d\nu, \quad \alpha \in \mathbb{R}. \quad 11$$

12 If  $P$  is not absolutely continuous with respect to  $Q$ , then  $\int \mathbb{I}\{p > 0, q = 0\} \, d\nu >$  12  
13  $0$ , and as a consequence,  $I_\alpha(P, Q) = \infty$  for  $\alpha \geq 1$ . A similar argument shows 13  
14 that  $I_\alpha(P, Q) = \infty$  if  $Q \ll P$  and  $\alpha \leq 0$ . Note that  $I_\alpha$  is well defined for  $\alpha = 0$  by 14  
15 taking the limit  $\alpha \rightarrow 0$  in the definition. Indeed, L'Hospital's rule implies that 15  
16

$$17 \quad \lim_{\alpha \rightarrow 0} I_\alpha(P, Q) = \int \log\left(\frac{p}{q}\right) q \, d\nu := K(P, Q) \quad 17$$

18 (with the above convention for the case where  $P \not\ll Q$ ), giving rise to the well- 18  
19 known Kullback–Leibler (KL) divergence measure from  $Q$  to  $P$ . The case 19  
20 with  $\alpha = 1$  corresponds to the KL divergence with the roles of  $P$  and  $Q$  20  
21 reversed. Note that the above definitions imply that the  $\alpha$ -divergence includes 21  
22 the Hellinger distance as a special case, in the sense that 22

$$23 \quad H^2(P, Q) = \frac{1}{2} I_{1/2}(P, Q). \quad 23$$

24 LEMMA 2.1: For probability measures  $P$  and  $Q$ , 24

$$25 \quad \max(\alpha, 1 - \alpha) I_\alpha(P, Q) \geq \frac{1}{2} I_{1/2}(P, Q) \quad 25$$

26 for every  $\alpha \in \mathbb{R}$ . 26

27 REMARK 2.1: Lemma 2.1 has some implications on a neighborhood system 27  
28 generated by the Hellinger distance. Consider the following neighborhood of 28  
29 a probability measure  $P$  whose radius in terms of  $I_\alpha$  is  $\delta > 0$ : 29

$$30 \quad B_{I_\alpha}(P, \delta) = \{Q \in \mathcal{M} : \sqrt{I_\alpha(Q, P)} \leq \delta\}. \quad 30$$

1 Lemma 2.1 implies that

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

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

$$6 C(L, U) = \left(\frac{1}{2} + L\right) \vee \left(\frac{1}{2} + U\right);$$

7 then by the above inequality,

$$8 (2.1) \quad \bigcup_{\alpha \in [\frac{1}{2} - L, \frac{1}{2} + U]} B_{I_{\alpha}}(P_0, \delta) \subset B_{I_{1/2}}(P_0, \sqrt{2C(L, U)}\delta).$$

9 That is, the union of the  $I_{\alpha}$ -based neighborhoods over  $\alpha \in [\frac{1}{2} - L, \frac{1}{2} + U]$  is  
10 covered by the Hellinger neighborhood  $B_{I_{1/2}}$  with a “margin” given by the mul-  
11 tiplicative constant  $\sqrt{2C(L, U)}$ . Equation (2.1) is important, since in what fol-  
12 lows we consider robustness of estimators against perturbation of  $P_0$  within  
13 its neighborhood, and it is desirable to use a neighborhood that is sufficiently  
14 large to accommodate a large class of perturbations. The inclusion relationship  
15 shows that the Hellinger-based neighborhood covers other neighborhood sys-  
16 tems based on  $I_{\alpha}$ ,  $\alpha \in [\frac{1}{2} - L, \frac{1}{2} + U]$  if the radii are chosen appropriately. It is  
17 easy to verify that (2.1) does not hold if the Hellinger distance  $I_{1/2}$  is replaced  
18 by  $I_{\alpha}$ ,  $\alpha \neq \frac{1}{2}$ , showing the special status of the Hellinger distance among the  
19  $\alpha$ -divergence family.

20 **REMARK 2.2:** Lemma 2.1 is a statement for every pair of measures  $(P, Q)$ ;  
21 thus it holds even if  $P \ll Q$  or  $Q \ll P$ . On the other hand, it is useful to con-  
22 sider the behavior of  $I_{\alpha}$  when one of the two measures is not absolutely con-  
23 tinuous with respect to the other. Consider a sequence of probability measures  
24  $\{P^{(n)}\}_{n \in \mathbb{N}}$ . 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 ev-  
25 ery  $\alpha' \in (0, 1)$ . But the reverse (i.e., reversing the roles of  $\alpha$  and  $\alpha'$ ) is not true.  
26 If  $P^{(n)}$ ,  $n \in \mathbb{N}$  are not absolutely continuous with respect to  $P_0$ ,  $I_{\alpha'}(P^{(n)}, P_0) = \infty$   
27 for every  $\alpha' \geq 1$  even if  $\rho_{\alpha}(P^{(n)}, P_0) \rightarrow 0$  for  $\alpha \in (0, 1)$  (and a similar argument  
28 holds for  $\alpha' \leq 0$ ). This shows that  $I_{\alpha}$ -based neighborhoods with  $\alpha \notin (0, 1)$  are  
29 too small: there are measures that are outside of  $B_{I_{\alpha}}(P_0, \delta)$ ,  $\alpha \notin (0, 1)$  no mat-  
30 ter how large  $\delta$  is, or how close they are to  $P_0$  in terms of, say, the Hellinger  
31 distance  $H$ .

1     REMARK 2.3: The inequality in Lemma 2.1 might be of interest on its own, 1  
2 as it generalizes many inequalities in the literature. For  $\alpha = 1$  or  $0$ , it cor- 2  
3 responds to the well-known inequality between the KL divergence and the 3  
4 Hellinger distance 4

$$5 \quad (2.2) \quad H(P, Q)^2 \leq K(P, Q); \quad 5$$

7 see, for example, Pollard (2002, p. 62). Another commonly used definition of 7  
8 divergence between probability measures is the  $\chi^2$  distance. It is given, if  $P \ll$  8  
9  $Q \ll \nu$ , by  $\chi^2(P, Q) = \int \frac{(p-q)^2}{q} d\nu$ , and it is shown that 9

$$10 \quad (2.3) \quad H(P, Q)^2 \leq \chi^2(P, Q) \quad 10$$

11 (Reiss (1989)). This is implied by Lemma 2.1 by letting  $\alpha = 2$ . Proposition 3.1 11  
12 in Zhang (2006) is closer to our result in terms of its generality; it shows that 12  
13  $\max(\alpha, 1 - \alpha)I_\alpha(P, Q) \geq \frac{1}{2}I_{1/2}(P, Q)$  holds for  $\alpha \in [0, 1]$ , which covers (2.2) but 13  
14 not (2.3)<sup>2</sup>. Lemma 2.1 shows that this type of inequality holds for all  $\alpha \in \mathbb{R}$ . 14  
15 16

17 Beran (1977), considering a parametric model, proposed MHDE that 17  
18 minimizes the Hellinger distance between a model-based probability measure 18  
19 (from the parametric family) and a nonparametric probability measure esti- 19  
20 mate. An application of the MHDE procedure to the moment condition model 20  
21 (1.1) yields a computationally simple procedure as follows. Let  $P_n$  denote the 21  
22 empirical measure of observations  $\{x_i\}_{i=1}^n$ .  $P_n$  is an appropriate model-free esti- 22  
23 mator in our construction of the MHDE. Let 23  
24 25

$$25 \quad \mathcal{P}_\theta = \left\{ P \in \mathcal{M} : \int g(x, \theta) dP = 0 \right\} \quad 25$$

26 and 26

$$27 \quad (2.4) \quad \mathcal{P} = \bigcup_{\theta \in \Theta} \mathcal{P}_\theta; \quad 27$$

28 then the MHDE, denoted by  $\hat{\theta}$ , is defined to be a parameter value that solves 28  
29 the optimization problem 29

$$30 \quad \inf_{\theta \in \Theta} \inf_{P \in \mathcal{P}_\theta} H(P, P_n) = \inf_{P \in \mathcal{P}} H(P, P_n). \quad 30$$

31 By convex duality theory (Kitamura (2006)), the objective function has the fol- 31  
32 lowing representation: 32  
33 34

$$35 \quad \inf_{P \in \mathcal{P}_\theta} H(P, P_n) = \max_{\gamma \in \mathbb{R}^m} -\frac{1}{n} \sum_{i=1}^n \frac{1}{1 + \gamma' g(x_i, \theta)}. \quad 35$$

36 <sup>2</sup>Zhang (2006) also derived a lower bound for the Hellinger distance in terms of  $I_\alpha$ . 36  
37 38  
39 40  
41 42  
43 44

1 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  
2 easy to compute.

3 It is easy to verify that we can obtain the MHDE as a Generalized Em-  
4 pirical Likelihood (GEL) estimator by letting  $\gamma = -1/2$  in equation (2.6)  
5 of Newey and Smith (2004). Asymptotic properties of the (G)EL estima-  
6 tors for  $\theta_0$  in (1.1), when data drawn from  $P_0$  are observed, are well under-  
7 stood (see, e.g., Kitamura and Stutzer (1997), Smith (1997), Imbens, Spady,  
8 and Johnson (1998), Newey and Smith (2004)). Let  $G = E_{P_0}[\partial g(x, \theta_0)/\partial \theta']$ ,  
9  $\Omega = E_{P_0}[g(x, \theta_0)g(x, \theta_0)']$ , and  $\Sigma = G'\Omega^{-1}G$ . Then

$$(2.5) \quad \sqrt{n}(\hat{\theta}_{\text{GEL}} - \theta_0) \xrightarrow{d} N(0, \Sigma^{-1}).$$

10  
11 It follows that the MHDE and other GEL estimators are semiparametrically  
12 efficient in the absence of data perturbation. At the same time, the MHDE  
13 possesses a distinct property of being asymptotic optimal robust if observations  
14 are drawn from a perturbed version of  $P_0$ , as we shall see in the next section.  
15  
16

### 17 3. ROBUST ESTIMATION THEORY

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

$$19 \quad T(P) = \arg \min_{\theta \in \Theta} \max_{\gamma \in \mathbb{R}^m} - \int \frac{1}{1 + \gamma'g(x, \theta)} dP;$$

20 then the MHDE can be interpreted as the value of functional  $T$  evaluated at  
21 the empirical measure  $P_n$ . In other words, each realization of  $P_n$  completely de-  
22 termines the value of the MHDE  $\hat{\theta}$ . To make the dependence explicit, we write  
23  $\hat{\theta} = T(P_n)$ , and study properties of the mapping  $T: \mathcal{M} \rightarrow \Theta$ . This definition of  
24  $T(\cdot)$ , however, causes a technical difficulty when the distribution of  $g(x, \theta)$  is  
25 unbounded for some  $\theta \in \Theta$  and  $P \in \mathcal{M}$ . To overcome this technical difficulty,  
26 we introduce the following mapping defined by a trimmed moment function:  
27  
28

$$29 \quad \bar{T}(Q) = \arg \min_{\theta \in \Theta} \inf_{P \in \bar{\mathcal{P}}_\theta, P \ll Q} H(P, Q),$$

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

$$32 \quad \bar{\mathcal{P}}_\theta = \left\{ P \in \mathcal{M} : \int g(x, \theta) \mathbb{I}\{x \in \mathcal{X}_n\} dP = 0 \right\},$$

$$33 \quad \mathcal{X}_n = \left\{ x \in \mathcal{X} : \sup_{\theta \in \Theta} |g(x, \theta)| \leq m_n \right\},$$

34 with the indicator function  $\mathbb{I}\{\cdot\}$  and the Euclidean norm  $|\cdot|$ ; that is,  $\mathcal{X}_n$  is a trim-  
35 ming set to bound the moment function and  $\bar{\mathcal{P}}_\theta$  is a set of probability measures  
36 satisfying the bounded moment condition  $E_P[g(x, \theta) \mathbb{I}\{x \in \mathcal{X}_n\}] = 0$ .  
37  
38  
39  
40  
41  
42  
43  
44

Let  $\tau: \Theta \rightarrow \mathbb{R}$  be a possibly nonlinear transformation of the parameter. We first focus on the estimation problem of the transformed scalar parameter  $\tau(\theta_0)$  and 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  $P_0$ ,

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

The transformation  $\tau$  to a scalar, as used by Rieder (1994), is convenient in calculating squared biases and MSEs. One may, for example, let  $\tau(\theta) = c'\theta$  using a constant  $p$ -vector  $c$ . Lemma A.1(ii) guarantees that, for each  $r > 0$ , the value  $\bar{T}(Q)$  exists for all  $Q \in B_H(P_0, r/\sqrt{n})$  and all  $n$  large enough.

ASSUMPTION 3.1: *Suppose the following conditions hold:*

- (i)  $\{x_i\}_{i=1}^n$  is independent and identically distributed (i.i.d.);
- (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. in  $\mathcal{N}$ ,  $\sup_{x \in \mathcal{X}_n, \theta \in \mathcal{N}} |\partial g(x, \theta)/\partial \theta| = o(n^{1/2})$ , 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  $n^{-1/2}m_n^{1+\varepsilon} = O(1)$  for some  $0 < \varepsilon < 2$  as  $n \rightarrow \infty$ ;
- (viii)  $\tau$  is continuously differentiable at  $\theta_0$ .

Assumption 3.1(i)–(vi) is standard in the literature of the 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 assumption is stronger than the one to derive the asymptotic distribution in (2.5). 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) or Newey (1990)) of  $\mathcal{P}$  in (2.4) such that  $P_{\theta_0, 0} = P_0$  and  $P_{\theta_0 + t/\sqrt{n}, \zeta_n} \in B_H(P_0, r/\sqrt{n})$  holds for  $\zeta_n = O(n^{-1/2})$  eventually.



1 (i)  $T_a$  is called *Fisher consistent* if, for every  $\{P_{\theta_n, \zeta_n}\}_{n \in \mathbb{N}}$  and  $t \in \mathbb{R}^p$ ,

$$2 \quad (3.1) \quad \sqrt{n}(T_a(P_{\theta_0+t/\sqrt{n}, \zeta_n}) - \theta_0) \rightarrow t. \quad 3$$

4 (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

$$5 \quad (3.2) \quad \sqrt{n}(T_a(P_n) - T_a(P_{\theta_n, \zeta_n})) \xrightarrow{d} M \quad \text{under } P_{\theta_n, \zeta_n}, \quad 6$$

7 where the measure  $M$  does not depend on the sequence  $(\theta'_n, \zeta'_n)'$ . 7

8 Both conditions are weak and satisfied by GMM, (G)EL, and other standard 8  
 9 estimators. For example, the mapping  $T_a$  for the continuous updating GMM 9  
 10 estimator (CUE) is given by 10

$$11 \quad T_{\text{CUE}}(P) = \underset{\theta \in \Theta}{\operatorname{argmin}} \left[ \int g(x, \theta) dP \right]' \left[ \int g(x, \theta) g(x, \theta) dP \right]^{-1} \quad 11$$

$$12 \quad \times \left[ \int g(x, \theta) dP \right], \quad 12$$

13 and, under Assumption 3.1,  $T_{\text{CUE}}(P_{\theta_0+t/\sqrt{n}, \zeta_n}) = \theta_0 + t/\sqrt{n}$  for large  $n$ . 13  
 14 CUE therefore trivially satisfies (3.1). The regularity condition (3.2) is stan- 14  
 15 dard in the literature of semiparametric efficiency; see, for example, Bickel 15  
 16 et al. (1993). 16

17 The following theorem shows the optimal robustness of the (trimmed) 17  
 18 MHDE in terms of its maximum bias. 18

19 **THEOREM 3.1:** *Suppose that Assumption 3.1 holds.* 19

20 (i) *For every  $T_a$  that is Fisher consistent,* 20

$$21 \quad \liminf_{n \rightarrow \infty} \sup_{Q \in B_H(P_0, r/\sqrt{n})} n(\tau \circ T_a(Q) - \tau(\theta_0))^2 \geq 4r^2 B^*, \quad 21$$

22 *for each  $r > 0$ , where  $B^* = (\frac{\partial \tau(\theta_0)}{\partial \theta})' \Sigma^{-1} (\frac{\partial \tau(\theta_0)}{\partial \theta})$ .* 22

23 (ii) *The mapping  $\bar{T}$  is Fisher consistent and satisfies* 23

$$24 \quad \lim_{n \rightarrow \infty} \sup_{Q \in B_H(P_0, r/\sqrt{n})} n(\tau \circ \bar{T}(Q) - \tau(\theta_0))^2 = 4r^2 B^*, \quad 24$$

25 *for each  $r > 0$ .* 25

26 **REMARK 3.1:** The above result is concerned with deterministic properties 26  
 27 of  $T_a$  and  $T$ .  $T_a(Q)$  and  $T(Q)$  can be regarded as the (probability) limit of the 27  
 28 estimators  $T_a(P_n)$  and  $T(P_n)$  under  $Q$ , and therefore the terms evaluated here 28  
 29 29  
 30 30  
 31 31  
 32 32  
 33 33  
 34 34  
 35 35  
 36 36  
 37 37  
 38 38  
 39 39  
 40 40  
 41 41  
 42 42  
 43 43  
 44 44

1 correspond to the bias of each estimator due to the deviation of  $Q$  from  $P_0$ . 1  
2 The theorem says that, in the class of all mappings that are Fisher consistent, 2  
3 the mapping  $\bar{T}$  has the smallest maximum bias over the set  $B_H(P_0, r/\sqrt{n})$ . The 3  
4 (trimmed version of) the Hellinger-based mapping  $\bar{T}$  is therefore optimally 4  
5 robust asymptotically in a minimax sense. The term  $4r^2B^*$  provides a sharp 5  
6 lower bound for maximum squared bias, and it is attained by  $\bar{T}$ . 6  
7

8 **REMARK 3.2:** The theorem is concerned with the trimmed version of the 8  
9 MHDE. It avoids the complications associated with the existence of  $T(Q)$  9  
10 for certain  $Q$ 's. If the support of  $\sup_{\theta \in \Theta} |g(x, \theta)|$  is bounded under every  $Q \in$  10  
11  $B_H(P_0, r/\sqrt{n})$  for large enough  $n$  (e.g., if the moment function  $g$  is bounded), 11  
12 then we do not need the trimming term  $\mathbb{I}\{x \in \mathcal{X}_n\}$ . In this case, the mapping  $T$  12  
13 without trimming has the above optimal robust property. 13  
14

15 **REMARK 3.3:** The index  $n$  in the statement of Theorem 3.1 simply parame- 15  
16 terizes how close  $Q \in B_H(P_0, r/\sqrt{n})$  and  $P_0$  are, and does not have to be inter- 16  
17 preted as the sample size. The next theorem, however, is concerned with MSEs 17  
18 and the index  $n$  represents the sample size there. 18  
19

20 The next theorem is our main result, which is concerned with (the supremum 20  
21 of) the MSE of the minimum Hellinger distance estimator  $\hat{\theta} = T(P_n)$  and other 21  
22 competing estimators. Let 22  
23

$$24 \quad (3.3) \quad \bar{B}_H(P_0, r/\sqrt{n}) = B_H(P_0, r/\sqrt{n}) \cap \left\{ Q \in \mathcal{M} : E_Q \left[ \sup_{\theta \in \Theta} |g(x, \theta)|^n \right] < \infty \right\}. \quad 24$$

25 We use the notation  $P^{\otimes n}$  to denote the  $n$ -fold product measure of a probability 25  
26 measure  $P$ . 26  
27

28 **THEOREM 3.2:** *Suppose that Assumption 3.1 holds.* 28  
29

30 (i) *For every Fisher consistent and regular mapping  $T_a$ ,* 30  
31

$$32 \quad \lim_{b \rightarrow \infty} \liminf_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n(\tau \circ T_a(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} \quad 32$$

$$33 \quad \geq (1 + 4r^2)B^*, \quad 33$$

$$34 \quad \geq (1 + 4r^2)B^*, \quad 34$$

$$35 \quad \geq (1 + 4r^2)B^*, \quad 35$$

$$36 \quad \geq (1 + 4r^2)B^*, \quad 36$$

37 *for each  $r > 0$ .* 37

38 (ii) *The mapping  $T$  is Fisher consistent and regular, and the MHDE  $\hat{\theta} = T(P_n)$*  38  
39 *satisfies* 39  
40

$$41 \quad \lim_{b \rightarrow \infty} \lim_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n(\tau \circ T(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} = (1 + 4r^2)B^*, \quad 41$$

$$42 \quad \lim_{b \rightarrow \infty} \lim_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n(\tau \circ T(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} = (1 + 4r^2)B^*, \quad 42$$

$$43 \quad \lim_{b \rightarrow \infty} \lim_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n(\tau \circ T(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} = (1 + 4r^2)B^*, \quad 43$$

$$44 \quad \lim_{b \rightarrow \infty} \lim_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge n(\tau \circ T(P_n) - \tau(\theta_0))^2 dQ^{\otimes n} = (1 + 4r^2)B^*, \quad 44$$

44 *for each  $r > 0$ .* 44

1     REMARK 3.4: This theorem establishes an asymptotic minimax optimal- 1  
2     ity property of the MHDE, in terms of MSE among all the estimators, 2  
3     that satisfies the two conditions in Definition 3.1. Note that the expression 3  
4      $\sup_{Q \in \tilde{B}_H(P_0, r/\sqrt{n})} \int b \wedge n(\tau \circ T_a(P_n) - \tau(\theta_0))^2 dQ^{\otimes n}$  is the maximum finite sam- 4  
5     ple MSE of  $T_a(P_n)$ . Thus our criterion for evaluating  $T_a$  (and  $T$ ) is the limit of 5  
6     its maximum finite sample MSE. Taking the supremum over  $B_H$  before letting 6  
7      $n$  go to infinity is important for capturing finite sample robustness properties. 7  
8     The method of calculating the truncated MSE first, then letting  $b \rightarrow \infty$ , is stan- 8  
9     dard in the literature of robust estimation, but is also used in general contexts; 9  
10    see, for example, [Bickel \(1981\)](#) and [LeCam and Yang \(1990\)](#). Once again, we 10  
11    are able to derive a sharp lower bound for the maximum MSE and show that 11  
12    it is achieved by the MHDE  $\hat{\theta} = T(P_n)$ . 12  
13

14     REMARK 3.5: Unlike in Theorem 3.1, optimality is achieved by the un- 14  
15     trimmed version of the MHDE. Note that  $T(P_n)$  exists for large  $n$  under As- 15  
16     sumption 3.1, in contrast to our discussion in Remark 3.2 on Theorem 3.1. 16  
17     Theorem 3.2, however, restricts the robustness neighborhood by an extra re- 17  
18     quirement as in (3.3). This is useful in showing that the untrimmed MHDE 18  
19     achieves the lower bound. 19  
20

21     REMARK 3.6: Theorem 3.2 proves that the MHDE is asymptotically opti- 21  
22     mally robust over a sequence of infinitesimal neighborhoods. Note that the 22  
23     Hellinger neighborhood over which the maximum of MSE is taken is nonpara- 23  
24     metric, in the sense that potential deviations from  $P_0$  cannot be indexed by a 24  
25     finite dimensional parameter. That is, our robustness concept demands uni- 25  
26     form robustness over a nonparametric, infinitesimal neighborhood. The use of 26  
27     infinitesimal neighborhoods, where the radius of the Hellinger ball shrinks at 27  
28     the rate  $n^{1/2}$ , is useful in balancing the magnitude of bias and variance in our 28  
29     asymptotics. If one uses a fixed, global neighborhood, then the bias term would 29  
30     dominate the behavior of estimators. This may fail to provide a good approxi- 30  
31     mation of finite sample behavior in actual applications, since in reality it would 31  
32     be reasonable to be concerned with both the stochastic fluctuation of estima- 32  
33     tors and their deterministic bias due to, say, data contamination. We note that 33  
34     there is a related but distinct literature on the asymptotics theory when the 34  
35     model is globally misspecified, as in [White \(1982\)](#), who considered paramet- 35  
36     ric MLE. [Kitamura \(1998, 2002\)](#) offered such analysis for conditional and un- 36  
37     conditional moment condition models. Moreover, [Schennach \(2007\)](#) provided 37  
38     novel and potentially very useful results of EL estimators and their variants in 38  
39     misspecified moment condition models. We regard our paper as a complement 39  
40     to, rather than a substitute for, the results obtained in these papers. There are 40  
41     fundamental differences between the characteristics of the problems the cur- 41  
42     rent paper considers and those of the papers on misspecification. First, our ob- 42  
43     ject of interest is  $\theta_0$ , not a pseudo-true value, as we consider data perturbation 43  
44     rather than model misspecification. Second, the nature of our analysis is local 44

1 and, therefore, the parameter value  $\theta_0$  in (1.1) is still identified asymptotically. 1  
2 Third, as noted above, we consider uniform robustness over a nonparametric 2  
3 neighborhood. The papers cited above consider pointwise problems. There- 3  
4 fore our approach deals with phenomena that are very different from the ones 4  
5 analyzed in the literature of misspecified models. 5  
6

7 **REMARK 3.7:** We have seen in Remark 2.1 that the Hellinger neighborhood 7  
8  $B_H$  has nice and distinct properties, in particular the inclusion relationship 8  
9 (2.1). The Hellinger neighborhood  $B_H$  is commonly used in the literature of ro- 9  
10 bust estimation (of parametric models); see, for example, [Beran \(1977\)](#), [Bickel](#) 10  
11 [\(1981\)](#), and [Rieder \(1994\)](#). We should note, however, that other neighborhood 11  
12 systems have been used in the literature as well. For example, one may replace 12  
13 the Hellinger distance  $H$  with the Kolmogorov–Smirnov (KS) distance in the 13  
14 definition of  $B_H$ . As [Beran \(1984\)](#) noted, however, to guarantee robustness in 14  
15 the Kolmogorov–Smirnov neighborhood system, one needs 15

16 “to use minimum distance estimates based on the Kolmogorov–Smirnov metric or a dis- 16  
17 tance weaker than the Kolmogorov–Smirnov metric . . . The general principle here is that 17  
18 the estimation distance be no stronger than the distance describing the contamination 18  
19 neighborhood. . . .” 19

20 [Donoho and Liu \(1988\)](#) developed a general theory of the above point. 20  
21 What this means is that an estimator that is robust against perturbations within 21  
22 Kolmogorov–Smirnov neighborhoods has to be minimizing the KS (or weaker) 22  
23 distance. The “minimum KS estimator” for the moment restriction model 23  
24 would be indeed robust, but it cannot be semiparametrically efficient when the 24  
25 model assumption holds. Therefore, unlike the moment restriction MHDE, 25  
26 the estimator is not “robust and efficient.” Another drawback is its computa- 26  
27 tion, since, unlike the moment restriction MHDE, no convenient algorithm to 27  
28 minimize the Kolmogorov–Smirnov distance under the moment restriction is 28  
29 known in the literature. It should be noted that the moment restriction MHDE 29  
30 is efficient in the sense that it achieves the semiparametric efficiency bound. It 30  
31 does not have the desirable higher order properties of EL ([Newey and Smith](#) 31  
32 [\(2004\)](#)) or the ETEL estimator proposed by [Schennach \(2007\)](#). 32  
33

34 The above MSE theorem conveniently summarizes the desirable robustness 34  
35 properties of the MHDE in terms of both (deterministic) bias and variance. It 35  
36 has, however, some limitations. First, its minimaxity result is obtained within 36  
37 Fisher consistent and regular estimators. While these requirements are weak, 37  
38 it might be of interest to expand the class of estimators. More importantly, 38  
39 implicit in the MSE-based analysis is that we are interested in  $L^2$ -loss. One may 39  
40 wish to use other types of loss functions, however, and it is of interest to see 40  
41 whether the above minimax results can be extended to a larger class of loss. The 41  
42 next theorem addresses these two issues. Of course, the MSE has an advantage 42  
43 of subsuming the bias and the variance in one measure. To deal with general 43  
44 loss functions, the next theorem focuses on the risk of estimators around a 44

1 Fisher consistent mapping evaluated at the perturbed measure  $Q$ . This can be 1  
2 regarded as calculating the risk of the first bracket of the decomposition (1.2), 2  
3 that is, the stochastic part of the deviation of the estimator from the parameter 3  
4 of interest  $\theta_0$ . 4

5 Let  $\mathcal{S}$  be a set of all estimators, that is, the set of all  $\bar{\mathbb{R}}^p$ -valued measur- 5  
6 able functions. We now investigate robust risk properties of this large class of 6  
7 estimators. The loss function we consider satisfies the following weak require- 7  
8 ments. 8

9  
10 ASSUMPTION 3.2: *The loss function  $\ell: \bar{\mathbb{R}}^p \rightarrow [0, \infty]$  is (i) symmetric subconvex 10  
11 (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); 11  
12 (ii) upper semicontinuous at infinity; and (iii) continuous on  $\bar{\mathbb{R}}^p$ . 12*

13  
14 We now present an optimal risk property for the MHDE. 14

15  
16 THEOREM 3.3: *Suppose that Assumptions 3.1 and 3.2 hold. 16*

17 (i) *For every Fisher consistent mapping  $T_a$ , 17*

$$18 \lim_{b \rightarrow \infty} \lim_{r \rightarrow \infty} \liminf_{n \rightarrow \infty} \inf_{S_n \in \mathcal{S}} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge \ell(\sqrt{n}(S_n - \tau \circ T_a(Q))) dQ^{\otimes n} 19$$

$$20 \geq \int \ell dN(0, B^*). 20$$

21  
22  
23  
24 (ii) *The mapping  $T$  is Fisher consistent and the MHDE  $\hat{\theta} = T(P_n)$  satisfies 24*

$$25 \lim_{b \rightarrow \infty} \lim_{n \rightarrow \infty} \sup_{Q \in \bar{B}_H(P_0, r/\sqrt{n})} \int b \wedge \ell(\sqrt{n}(\tau \circ T(P_n) - \tau \circ \bar{T}(Q))) dQ^{\otimes n} 26$$

$$27 = \int \ell dN(0, B^*), 27$$

28  
29  
30  
31  
32 for all  $r > 0$ . 32

33  
34 Note that Theorem 3.3(ii) remains valid if  $T(P_n)$  is replaced by  $\bar{T}(P_n)$ . This 34  
35 theorem shows that the MHDE is once again optimally robust even for the general 35  
36 risk criterion, and this holds in the class of essentially all possible estimators. 36  
37 As noted above, the result is concerned with the stochastic component of the 37  
38 decomposition (1.2). Theorem 3.1 has already established that the MHDE 38  
39 is optimal in terms of its bias, that is, the deterministic part of the decompo- 39  
40 sition (1.2) in the second bracket. The latter result does not depend on a spe- 40  
41 cific loss function. Thus the MHDE enjoys general optimal robust properties 41  
42 under a quite general setting, in terms of both the stochastic component and 42  
43 the deterministic component. Note that analyzing these two parts separately is 43  
44 common in the literature of robust statistics: see, for example, Rieder (1994). 44

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.

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  $g$  when the original DGP is perturbed 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  $\delta = 0.25\sin(2\pi\omega)$  and let  $\omega$  vary over  $\omega_j = j/64, j = 0, \dots, 63$ . This yields 64 different designs; for each of them, 5000 replications are performed and RMSE and  $\Pr\{|\hat{\theta} - \theta_0| > 0.5\}$  are calculated. We consider the following estimators: empirical likelihood (EL), MHDE, exponential tilting (ET), GMM (GMM2), and continuously updated GMM (CUE). GMM2 is calculated following the standard two step procedure where the initial estimate is obtained from identity weighting. CUE’s performance is extremely sensitive to data perturbations considered here; its RMSE is much higher than that of the other estimators. For convenience, we only plot the results for EL, MHDE, ET, and GMM2 in displaying their RMSEs. The results are presented in Figure 1. In the left panel, each curve represents the RMSE of a particular estimator as a function of  $\omega_j$ . The right panel (labeled “Pr”) displays the simulated probability of an estimator deviating from the target  $\theta_0 = 3$  by more than 0.5.

While RMSE is a potentially informative measure, it can be highly misleading, as some of the estimators may not have finite moments. We thus focus on the results for deviation probabilities. The performance of CUE clearly indicates its lack of robustness against data perturbations. We also see that GMM2 is affected by perturbations much more than EL, MHDE, and ET, except for

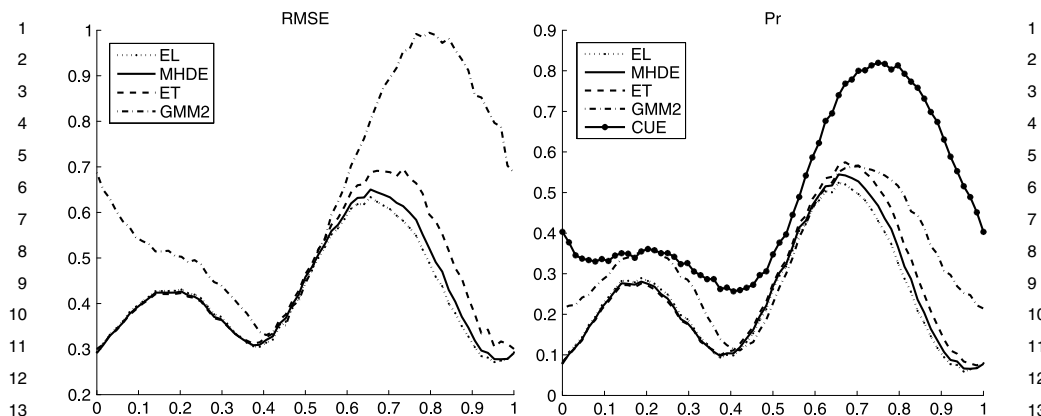


FIGURE 1.—Local neighborhood of the true model. “Pr” denotes  $\Pr\{|\hat{\theta} - \theta_0| > 0.5\}$ .

the values of  $\omega$ 's between 0.4 and 0.6, where the performances of the estimators other than CUE are rather close. ET seems to perform a little worse than MHDE and EL.

One needs to 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 that 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 exhibit more instability throughout the different perturbation designs. 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 nonrobust, has  $\gamma = 1$ , a value that is much higher than the optimally robust  $\gamma = -1/2$  of the MHDE.

### 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, and therefore imposes mild assumptions. Yet it is reasonable to expect that the probability law of observations may have some deviations from the ideal

1 distribution as modeled by the moment restriction model. It is then sensible 1  
2 to seek estimation procedures that are robust against slight perturbations in 2  
3 the probability measure that generates observations, which can be caused by, 3  
4 for example, data contamination. Our main theoretical result shows that the 4  
5 minimum Hellinger distance estimator (MHDE) possesses optimal minimax 5  
6 robust properties. Moreover, it remains semiparametrically efficient when the 6  
7 model assumptions hold. Convenient numerical algorithms for its implementa- 7  
8 tion are provided. Our simulation results indicate that GMM can be highly 8  
9 sensitive to data perturbations. The performance of the MHDE remains stable 9  
10 over a wide range of simulation designs, which is in accordance with our 10  
11 theoretical findings. 11

12 The results obtained in this paper are concerned with estimation, though 12  
13 it might be potentially possible to extend our robustness theory to parameter 13  
14 testing problems. It is of practical importance to consider robust methods 14  
15 for testing and confidence interval calculations so that the results of statisti- 15  
16 cal inference for moment restriction models are reliable and not too sensitive 16  
17 to departures from model assumptions. Interestingly, there exists a literature 17  
18 on parametric robust inference based on the MHDE method. We plan to investi- 18  
19 gate robust testing procedure in moment condition models in our future 19  
20 research. 20  
21

## 22 REFERENCES 22

- 23  
24 BERAN, R. (1977): “Minimum Hellinger Distance Estimates for Parametric Models,” *The Annals*  
25 *of Statistics*, 5, 445–463. MR0448700[2,6,12,15] 25 <LS\_link>  
26 ——— (1984): “Minimum Distance Procedures,” in *Handbook of Statistics*, ed. by P. Krishnaiah  
27 and P. Sen. Amsterdam: Elsevier Science, 741–754. MR0831734[12] 26  
28 BICKEL, P. J. (1981): “Quelques aspects de la statistique robuste”, in *Ecole d’Eté de Probabilités*  
29 *de Saint Flour IX 1979*, ed. by P. Hennequin. Berlin: Springer, 1–72. MR0637470[11,12] 28  
30 BICKEL, P., C. KLASSEN, Y. RITOV, AND J. WELLNER (1993): *Efficient and Adaptive Estimation for*  
31 *Semiparametric Models*. Baltimore: Johns Hopkins Press. MR1245941[8,9] 30  
32 DONOHO, D., AND R. LIU (1988): “The “Automatic” Robustness of Minimum Distance Func-  
33 tional,” *The Annals of Statistics*, 16, 552–586. MR0947562[12] 31 <LS\_link>  
34 HALL, P., AND J. L. HOROWITZ (1996): “Bootstrap Critical Values for Tests Based on Generalized  
35 Method of Moments Estimators,” *Econometrica*, 64, 891–916. MR1399222[14] 32 <LS\_link>  
36 HANSEN, L. P. (1982): “Large Sample Properties of Generalized Methods of Moments Estima-  
37 tors,” *Econometrica*, 50, 1029–1054. MR0666123[3] 33 <LS\_link>  
38 IMBENS, G. W., R. H. SPADY, AND P. JOHNSON (1998): “Information Theoretic Approaches to  
39 Inference in Moment Condition Models,” *Econometrica*, 66, 333–357. MR1612246[7] 34 <LS\_link>  
40 KITAMURA, Y. (1998): “Comparing Misspecified Dynamic Econometric Models Using Nonpara-  
41 metric Likelihood,” Working Paper, Department of Economics, University of Wisconsin. [11] 35 <LS\_link>  
42 ——— (2002): “A Likelihood-Based Approach to the Analysis of a Class of Nested and Non-  
43 Nested Models,” Working Paper, Department of Economics, University of Pennsylvania. [11] 36 <LS\_link>  
44 ——— (2006): “Empirical Likelihood Methods in Econometrics: Theory and Practice,” in *Ad-*  
45 *vances in Economics and Econometrics: Theory and Applications, Ninth World Congress*, ed. by  
46 R. Blundell, W. K. Newey, and T. Persson. Cambridge: Cambridge University Press. [6] 37  
47 KITAMURA, Y., AND M. STUTZER (1997): “An Information Theoretic Alternative to Generalized  
48 Method of Moments Estimation,” *Econometrica*, 65 (4), 861–874. MR1458431[7] 38 <LS\_link>



- <uncited> 1 KITAMURA, Y., T. OTSU, AND K. EVDOKIMOV (2013): “Supplement to ‘Robustness, 1  
2 Infinitesimal Neighborhoods, and Moment Restrictions,’” *Econometrica Supplemental 2  
3 Material*, 81, [http://www.econometricsociety.org/ecta/Supmat/8617\\_proofs.pdf](http://www.econometricsociety.org/ecta/Supmat/8617_proofs.pdf); [http://www.econometricsociety.org/ecta/Supmat/8617\\_data\\_and\\_programs.zip](http://www.econometricsociety.org/ecta/Supmat/8617_data_and_programs.zip). 3  
4  
5 LECAM, L., AND G. YANG (1990): *Asymptotics in Statistics: Some Basic Concepts*. New York: 4  
5 Springer. MR1066869[11] 5  
6 NEWEY, W. K. (1990): “Semiparametric Efficiency Bounds,” *Journal of Applied Econometrics*, 5, 6 <LS\_link>  
7 99–135. [8] 7  
8 NEWEY, W. K., AND R. J. SMITH (2004): “Higher Order Properties of GMM and Generalized 8  
9 Empirical Likelihood Estimators,” *Econometrica*, 72, 219–255. MR2031017[7,12,15] 9 <LS\_link>  
10 POLLARD, D. (2002): *A User’s Guide to Measure Theoretic Probability*. Cambridge: Cambridge 10  
11 University Press. [6] 11  
12 REISS, R.-D. (1989): *Approximate Distributions of Order Statistics*. New York: Springer-Verlag. 12  
13 MR0988164[6] 13  
14 RIEDER, H. (1994): *Robust Asymptotic Statistics*. New York: Springer-Verlag. MR1284041[2,8,12, 14  
15 13] 15  
16 SCHENNACH, S. M. (2007): “Point Estimation With Exponentially Tilted Empirical Likelihood,” 16 <LS\_link>  
17 *The Annals of Statistics*, 35, 634–672. MR2336862[11,12] 17  
18 SMITH, R. J. (1997): “Alternative Semi-Parametric Likelihood Approaches to Generalized 18 <LS\_link>  
19 Method of Moments Estimation,” *Economic Journal*, 107, 503–519. [7] 19 <LS\_link>  
20 WHITE, H. (1982): “Maximum Likelihood Estimation of Misspecified Models,” *Econometrica*, 20  
21 50, 1–25. MR0640163[11] 21  
22 ZHANG, T. (2006): “From  $\varepsilon$ -Entropy to KL-Entropy: Analysis of Minimum Information Com- 22  
23 plexity Density Estimation,” *The Annals of Statistics*, 34, 2180–2210. MR2291497[6] 23 <LS\_link>  
24  
25 *Cowles Foundation for Research in Economics, Yale University, New Haven,* 24  
26 *CT 06520, U.S.A.; [yuichi.kitamura@yale.edu](mailto:yuichi.kitamura@yale.edu),* 25  
27 *Cowles Foundation for Research in Economics, Yale University, New Haven,* 26  
28 *CT 06520, U.S.A.; [taisuke.otsu@yale.edu](mailto:taisuke.otsu@yale.edu),* 27  
29 *and* 28  
30 *Dept. of Economics, Princeton University, Princeton, NJ 08544, U.S.A.;* 29  
31 *[kevdokim@princeton.edu](mailto:kevdokim@princeton.edu).* 30  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44

Manuscript received June, 2009; final revision received August, 2012.

## 1 THE LIST OF SOURCE ENTRIES RETRIEVED FROM MATHSCINET 1

2 **The list of entries below corresponds to the Reference section of your article and was retrieved**  
3 **from MathSciNet applying an automated procedure. Please check the list and cross out those**  
4 **entries which lead to mistaken sources. Please update your references entries with the data from**  
5 **the corresponding sources, when applicable. More information can be found in the [support page](#).**

- 6  
7 BERAN, R. (1977). Minimum Hellinger distance estimates for parametric models. *Ann. Statist.* **5**  
8 445–463. [MR0448700](#)  
9 BERAN, R. (1984). Minimum distance procedures. In *Nonparametric methods*. Handbook of  
10 Statist., Vol. 4. North-Holland, Amsterdam, 741–754. [MR0831734](#)  
11 BICKEL, P. J., KLAASSEN, C. A. J., RITOV, Y., AND WELLNER, J. A. (1993). *Efficient and adaptive*  
12 *estimation for semiparametric models*. Johns Hopkins Series in the Mathematical Sciences.  
13 Johns Hopkins University Press, Baltimore, MD. [MR1245941](#)  
14 BICKEL, P. J. (1981). Quelques aspects de la statistique robuste. In *Ninth Saint Flour Probability*  
15 *Summer School—1979 (Saint Flour, 1979)*. Lecture Notes in Math., Vol. **876**. Springer, Berlin,  
16 1–72. [MR0637470](#)  
17 DONOHO, D. L. AND LIU, R. C. (1988). The “automatic” robustness of minimum distance func-  
18 tionals. *Ann. Statist.* **16** 552–586. [MR0947562](#)  
19 HALL, P. AND HOROWITZ, J. L. (1996). Bootstrap critical values for tests based on generalized-  
20 method-of-moments estimators. *Econometrica* **64** 891–916. [MR1399222](#)  
21 HANSEN, L. P. (1982). Large sample properties of generalized method of moments estimators.  
22 *Econometrica* **50** 1029–1054. [MR0666123](#)  
23 IMBENS, G. W., SPADY, R. H., AND JOHNSON, P. (1998). Information-theoretic approaches to  
24 inference in moment condition models. *Econometrica* **66** 333–357. [MR1612246](#)  
25 Not Found!  
26 Not Found!  
27 (2006). *Advances in economics and econometrics: theory and applications. Vol. I*. Econometric  
28 Society Monographs, Vol. **41**. Cambridge University Press, Cambridge. Edited by Richard  
29 Blundell, Whitney K. Newey and Torsten Persson. [MR2352812](#)  
30 KITAMURA, Y. AND STUTZER, M. (1997). An information-theoretic alternative to generalized  
31 method of moments estimation. *Econometrica* **65** 861–874. [MR1458431](#)  
32 LE CAM, L. AND YANG, G. L. (1990). *Asymptotics in statistics*. Springer Series in Statistics.  
33 Springer-Verlag, New York. Some basic concepts. [MR1066869](#)  
34 Not Found!  
35 NEWAY, W. K. AND SMITH, R. J. (2004). Higher order properties of GMM and generalized  
36 empirical likelihood estimators. *Econometrica* **72** 219–255. [MR2031017](#)  
37 Not Found!  
38 REISS, R.-D. (1989). *Approximate distributions of order statistics*. Springer Series in Statistics.  
39 Springer-Verlag, New York. With applications to nonparametric statistics. [MR0988164](#)  
40 RIEDER, H. (1994). *Robust asymptotic statistics*. Springer Series in Statistics. Springer-Verlag,  
41 New York. [MR1284041](#)  
42 SCHENNACH, S. M. (2007). Point estimation with exponentially tilted empirical likelihood. *Ann.*  
43 *Statist.* **35** 634–672. [MR2336862](#)  
44 Not Found!  
45 WHITE, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica* **50**  
46 1–25. [MR0640163](#)  
47 ZHANG, T. (2006). From  $\epsilon$ -entropy to KL-entropy: analysis of minimum information complexity  
48 density estimation. *Ann. Statist.* **34** 2180–2210. [MR2291497](#)  
49  
50  
51  
52

1 META DATA IN THE PDF FILE 1

2 **Following information will be included as pdf file Document Properties:** 2

3 3

4 **Title** : Robustness, Infinitesimal Neighborhoods, and Moment Restric- 4  
tions

5 **Author** : Yuichi Kitamura, Taisuke Otsu, Kirill Evdokimov 5

6 **Subject** : Econometrica, Vol.0, No.00, ????, 0, 1-17 6

7 **Keywords**: Asymptotic Minimax Theorem, Hellinger distance, semiparamet- 7  
ric efficiency 8

9 9

10 THE LIST OF URI ADRESSES 10

11 11

12 12

13 **Listed below are all uri addresses found in your paper. The non-active uri addresses, if any, are** 13  
**indicated as ERROR. Please check and update the list where necessary. The e-mail addresses** 14  
**are not checked – they are listed just for your information. More information can be found in the** 15  
**support page.** 16

17 17

18 200 <http://www.econometricsociety.org/> [6:pp.0,0,1,1,1,1] OK 18

19 200 [http://www.econometricsociety.org/ecta/Supmat/8617\\_proofs.pdf](http://www.econometricsociety.org/ecta/Supmat/8617_proofs.pdf) [2:pp.18,18] OK 19

20 200 [http://www.econometricsociety.org/ecta/Supmat/8617\\_data\\_and\\_programs.zip](http://www.econometricsociety.org/ecta/Supmat/8617_data_and_programs.zip) [2:pp.18,18] OK 20

21 --- <mailto:yuichi.kitamura@yale.edu> [2:pp.19,19] Check skip 21

22 --- <mailto:taisuke.otsu@yale.edu> [2:pp.19,19] Check skip 22

23 --- <mailto:kevokim@princeton.edu> [2:pp.19,19] Check skip 23

24 24

25 25

26 26

27 27

28 28

29 29

30 30

31 31

32 32

33 33

34 34

35 35

36 36

37 37

38 38

39 39

40 40

41 41

42 42

43 43

44 44