

# $\sqrt{n}$ -consistent parameter estimation for systems of ordinary differential equations : bypassing numerical integration via smoothing

**Citation for published version (APA):**

Gugushvili, S., & Klaassen, C. A. J. (2010).  $\sqrt{n}$ -consistent parameter estimation for systems of ordinary differential equations : bypassing numerical integration via smoothing. (Report Eurandom; Vol. 2010033). Eurandom.

**Document status and date:**

Published: 01/01/2010

**Document Version:**

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

**Please check the document version of this publication:**

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

**General rights**

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

[www.tue.nl/taverne](http://www.tue.nl/taverne)

**Take down policy**

If you believe that this document breaches copyright please contact us at:

[openaccess@tue.nl](mailto:openaccess@tue.nl)

providing details and we will investigate your claim.

EURANDOM PREPRINT SERIES  
2010-033

**$\sqrt{n}$ -CONSISTENT PARAMETER ESTIMATION FOR SYSTEMS  
OF ORDINARY DIFFERENTIAL EQUATIONS:  
BYPASSING NUMERICAL INTEGRATION VIA SMOOTHING**

S. Gugushvili, C.A.J. Klaassen  
ISSN 1389-2355

**$\sqrt{n}$ -CONSISTENT PARAMETER ESTIMATION FOR SYSTEMS  
OF ORDINARY DIFFERENTIAL EQUATIONS: BYPASSING  
NUMERICAL INTEGRATION VIA SMOOTHING**

SHOTA GUGUSHVILI AND CHRIS A.J. KLAASSEN

ABSTRACT. We consider the problem of parameter estimation for a system of ordinary differential equations from noisy observations on a solution of the system. In case the system is nonlinear, as it typically is in practical applications, an analytic solution to it usually does not exist. Consequently, straightforward estimation methods like the ordinary least squares method depend on repetitive use of numerical integration in order to determine the solution of the system for each of the parameter values considered, and to find subsequently the parameter estimate that minimises the objective function. This induces a huge computational load to such estimation methods. We propose an estimator that is defined as a minimiser of an appropriate distance between a nonparametrically estimated derivative of the solution and the right-hand side of the system applied to a nonparametrically estimated solution. Our estimator bypasses numerical integration altogether and reduces the amount of computational time drastically compared to ordinary least squares. Moreover, we show that under suitable regularity conditions this estimation procedure leads to a  $\sqrt{n}$ -consistent estimator of the parameter of interest.

1. BRIEF INTRODUCTION

Many dynamical systems in science and applications are modelled by a  $d$ -dimensional system of ordinary differential equations, denoted as

$$(1) \quad \begin{cases} x'(t) = F(x(t), \theta), & t \in [0, 1], \\ x(0) = \xi, \end{cases}$$

where  $\theta$  is the unknown parameter of interest and  $\xi$  is the initial condition. With  $x_\theta(t)$  the solution vector corresponding to the parameter value  $\theta$ , we observe

$$Y_{ij} = x_{\theta_j}(t_i) + \epsilon_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, d,$$

where the observation times  $0 \leq t_1 < \dots < t_n \leq 1$  are known and the random variables  $\epsilon_{ij}$  have mean 0 and model measurement errors combined with latent random deviations from the idealised model (1). Under regularity conditions the

---

*Date:* July 23, 2010.

*2000 Mathematics Subject Classification.* Primary: 62F12, Secondary: 62G08, 62G20.

*Key words and phrases.* M-estimator;  $\sqrt{n}$ -consistency; Nonparametric regression; ODE system; Priestley-Chao estimator.

Most of the research was done while the first author was at EURANDOM, Eindhoven, The Netherlands.

ordinary least squares estimator

$$(2) \quad \tilde{\theta}_n = \operatorname{argmin}_{\eta} \sum_{i=1}^n \sum_{j=1}^d (Y_{ij} - x_{\eta j}(t_i))^2$$

of  $\theta$  is  $\sqrt{n}$ -consistent, at least theoretically. For systems (1) that do not have explicit solutions, one typically uses iterative procedures to approximate this ordinary least squares estimator. However, since every iteration in such a procedure involves numerical integration of the system (1) and since the number of iterations is typically very large, in practice it is often extremely difficult if not impossible to compute (2). Here we present a feasible and computationally much faster method to estimate the parameter  $\theta$ . To define our estimator we first construct kernel estimators

$$\hat{x}_j(t) = \sum_{i=1}^n (t_i - t_{i-1}) \frac{1}{b} K\left(\frac{t - t_i}{b}\right) Y_{ij}$$

of  $x_{\theta j}$  with  $K(\cdot)$  a kernel function and  $b = b_n$  a bandwidth. Now, our estimator  $\hat{\theta}_n$  of  $\theta$  is defined as

$$\hat{\theta}_n = \operatorname{argmin}_{\eta} \int_0^1 \|\hat{x}'(t) - F(\hat{x}(t), \eta)\|^2 w(t) dt,$$

where  $\|\cdot\|$  denotes the usual Euclidean norm and  $w(\cdot)$  is a weight function.

The main result of this paper is that this estimator is  $\sqrt{n}$ -consistent under mild regularity conditions. So, our estimator is comparable to the ordinary least squares estimator in statistical performance, but it avoids the computationally costly repeated use of numerical integration of (1).

## 2. INTRODUCTION

Let us introduce the contents of this paper in more detail. Systems of ordinary differential equations play a fundamental role in many branches of natural sciences, e.g. mathematical biology, see Edelman-Keshet (2005), biochemistry, see Voit (2000), or the theory of chemical reaction networks in general, see for instance Feinberg (1979) and Sontag (2001). Such systems usually depend on parameters, which in practice are often only approximately known, or are plainly unknown. Knowledge of these parameters is critical for the study of the dynamical system or process that the system of ordinary differential equations describes. Since these parameters usually cannot be measured directly, they have to be inferred from, as a rule, noisy measurements of various quantities associated with the process under study. More formally, in this paper we consider the following setting: let

$$(3) \quad \begin{cases} x'(t) = F(x(t), \theta), & t \in [0, 1], \\ x(0) = \xi, \end{cases}$$

be a system of autonomous differential equations depending on a vector of real-valued parameters. Here  $x(t) = (x_1(t), \dots, x_d(t))^T$  is a  $d$ -dimensional state variable,  $\theta = (\theta_1, \dots, \theta_p)^T$  denotes a  $p$ -dimensional parameter, while the column  $d$ -vector  $x(0) = \xi$  defines the initial condition. Whether the latter is known or unknown, is not relevant in the present context, as long as it stays fixed. Denote a solution to

(3) corresponding to parameter value  $\theta$  by  $x_\theta(t) = (x_{\theta 1}(t), \dots, x_{\theta d}(t))^T$ . Suppose that at known time instances  $0 \leq t_1 < \dots < t_n \leq 1$  noisy observations

$$(4) \quad Y_{ij} = x_{\theta j}(t_i) + \epsilon_{ij}, \quad i = 1, \dots, n, j = 1, \dots, d,$$

on the solution  $x_\theta$  are available. The random variables  $\epsilon_{ij}$  model measurement errors, but they might also contain latent random deviations from the idealized model (1). Such random deviations are often seen in real-world applications. Based on these observations, the goal is to infer the value of  $\theta$ , the parameter of interest.

A standard approach to estimation of  $\theta$  is based on the least squares method, see e.g. Hemker (1972) and Stortelder (1996). Assuming for simplicity  $d = 1$  for the moment, the least squares estimator is defined as a minimiser of the sum of squares, i.e.

$$\tilde{\theta}_n = \operatorname{argmin}_\eta R_n(\eta) = \operatorname{argmin}_\eta \sum_{i=1}^n (Y_i - x_\eta(t_i))^2.$$

If the measurement errors are Gaussian, then  $\tilde{\theta}_n$  coincides with the maximum likelihood estimator and is asymptotically efficient. Since the differential equations setting is covered by the general theory of nonlinear least squares, theoretical results available for the latter apply also in the differential equations setting and we refer e.g. to Jennrich (1969) and Wu (1981), or more generally to van de Geer (1990), van de Geer and Wegkamp (1996), and Pollard and Radchenko (2006) for a thorough treatment of the asymptotics of the nonlinear least squares estimator. Despite its appealing theoretical properties, in practice the performance of the least squares method can dramatically degrade if (3) is a nonlinear high-dimensional system and if  $\theta$  is high-dimensional. In such a case we have to face a nonlinear optimisation problem (quite often with many local minima) and search for a global minimum of the least squares criterion function  $R_n$  in a high-dimensional parameter space. The search process is most often done via gradient-based methods, e.g. the Levenberg-Marquardt method, see Levenberg (1963), or via random search algorithms, see Section 4.5.2 in Voit (2000) for a literature overview. Since nonlinear systems in general do not have solutions in closed form, use of numerical integration within a gradient-based search method and serious computational time associated with it seem to be inevitable. For instance, a relatively simple example of a four-dimensional system considered in Appendix 1 of Voit and Almeida (2004) demonstrates that the need to repeat numerical integration multiple times might increase the computational time for numerical integration up to 95% of the total computational time required for a gradient based optimisation method. Likewise, random search algorithms are also very costly computationally. The problems become aggravated for systems of ordinary differential equations that exhibit stiff behaviour, i.e. systems that are difficult to integrate via explicit numerical integration schemes, see e.g. Hairer and Wanner (1996) for a comprehensive treatment of methods of solving numerically stiff systems. Even if a system is not stiff for the true parameter value, during the numerical optimisation procedure one might pass the vicinity of parameters for which the system is stiff, which will necessarily slow down the optimisation process.

The Bayesian approach to estimation of  $\theta$ , see e.g. Gelman et al. (1996) and Girolami (2008), encounters similar huge computational problems. In the Bayesian approach one puts a prior on the parameter  $\theta$  and then obtains the posterior via Bayes' formula. The posterior contains all the information required in the Bayesian

paradigm and can be used to compute e.g. point estimates of  $\theta$  or Bayesian credible intervals. If  $\theta$  is high-dimensional, the posterior will typically not be manageable by numerical integration and one will have to resort to Markov Chain Monte Carlo (MCMC) methods. However, sampling from the posterior distribution for  $\theta$  via MCMC necessitates at each step numerical integration of the system (3), in case the latter does not have a closed form solution. Computational time might thus become a problem in this case as well. Also, since in general the likelihood surface will have a complex shape with many local optima, ripples, and ridges, see e.g. Girolami (2008) for an example, serious convergence problems might arise for MCMC samplers.

Yet another point is that in practice both the least squares method and the Bayesian approach require good initial guesses of the parameter values. If these are not available, then both approaches might have problems with convergence to the true parameter value within a reasonable amount of time. More generally, computational time will typically be a problem for any optimisation algorithm that relies on numerical integration of any relatively realistic nonlinear system of ordinary differential equations. One example is furnished by Kikuchi et al. (2003), where a system that consists of five differential equations and contains sixty parameters and that describes a simple gene regulatory network from Hlavacek and Savageau (1996) is considered. The optimisation algorithm (a genetic algorithm) was run for seven loops each lasting for about ten hours on the AIST CBRC Magi Cluster with 1040 CPUs (Pentium III 933 MHz)<sup>1</sup>. This amounted to a total of ca. 70,000 CPU hours. The authors also remarked that the gradient-based search algorithm would not be feasible in their setting at all.

A general overview of typical difficulties in parameter estimation for systems of ordinary differential equations is given in Ramsay et al. (2007), to which we refer for more details. For a recent overview of typical approaches to parameter estimation for systems of ordinary differential equations in biochemistry and associated challenges see e.g. Chou and Voit (2009).

To evade difficulties associated with the least squares method, or more precisely with numerical integration that it usually requires, a two-step method was proposed in Varah (1982). In the first step the solution  $x_\theta$  of (3) is estimated via considering estimation of the individual components  $x_{\theta_1}, \dots, x_{\theta_d}$  as nonparametric regression problems and using the regression spline method for estimation of these components. The derivatives of  $x_{\theta_1}, \dots, x_{\theta_d}$  are also estimated from the data by differentiating the estimators of  $x_{\theta_1}, \dots, x_{\theta_d}$  with respect to time  $t$ . Thus no numerical integration of the system (3) is needed. In the second step the obtained estimate of  $x_\theta$  and its derivative  $x'_\theta$  are plugged into (3) and an estimator of  $\theta$  is defined as a minimiser in  $\theta$  of an appropriate distance between the estimated left- and righthand sides of (3). Such an estimator of  $\theta$  is an M-estimator, see e.g. the classical monograph Huber (1981), or Chapter 7 of Bickel et al. (1998), Chapter 5 of van der Vaart (1998), and Chapter 3.2 of Wellner and van der Vaart (1996) for a more modern exposition of the theory of M-estimators. For a related approach to estimation of  $\theta$  see also Voit and Savageau (1982), as well as Voit and Almeida (2004), where a practical implementation based on neural networks is studied. The intuitive idea behind the use of this two-step estimator is clear: among all functions defined on  $[0, 1]$ , any reasonably defined distance between the left- and righthand

---

<sup>1</sup>See <http://www.cbrc.jp/magi> for the cluster specifications.

side of (3) is minimal (namely, it is zero) for the solution  $x_\theta$  of (3) and the true parameter value  $\theta$ . For estimates close enough in an appropriate sense to the solution  $x_\theta$ , the minimisation procedure will produce a minimiser close to the true parameter value, provided certain identifiability and continuity conditions hold. This intuitive idea was exploited in Brunel (2008), where a more general setting than the one in Varah (1982) was considered. Another paper in the same spirit as Varah (1982) is Liang and Wu (2008).

This two-step approach will typically lead to considerable savings in computational time, as unlike the straightforward least squares estimator, in its first step it just requires finding nonparametric estimates of  $x_\theta$  and  $x'_\theta$ , for which fast and numerically reliable recipes are available, whereas the gradient-based least squares method will still rely on successive numerical integrations of (3) for different parameter values  $\theta$  in order to find a global minimiser minimising the least squares criterion function. We refer to Voit and Almeida (2004) for a particular example demonstrating gains in the computational time achieved by the two-step estimator in comparison to the ordinary least squares estimator. When the righthand side  $F$  of (3) is linear in  $\theta_1, \dots, \theta_p$ , further simplifications will occur in the second step of the two-step estimation procedure, as one will essentially only have to face a weighted linear regression problem then. This is unlike the least squares approach, which cannot exploit linearity of  $F$  in  $\theta_1, \dots, \theta_p$ . However, we would also like to stress the fact that the two-step estimator does not necessarily have to be considered a competitor of either the least squares or the Bayesian approach. Indeed, since in practice both of these approaches require good initial guesses for parameter values, these can be supplied by the two-step estimator. In this sense the proposed two-step estimation approach can be thought of as complementing both the least squares and the Bayesian approaches. Moreover, an additional modified Newton-Raphson step suffices to arrive at an estimator that is asymptotically equivalent to the exact ordinary least squares estimator, as will be shown elsewhere.

Our exposition in the present paper is similar to that in Brunel (2008) to some degree, one of the differences being that instead of spline estimators we use kernel-type estimators for estimation of  $x_\theta$  and  $x'_\theta$ .<sup>2</sup> The conditions are also somewhat different. We hope that our contribution will motivate further research into the interesting topic of parameter estimation for systems of ordinary differential equations. There exists an alternative approach to the ones described here, which also employs nonparametric smoothing, see Ramsay et al. (2007). For information on its asymptotic properties we refer to Qi and Zhao (2010). For nonlinear systems this approach will typically reduce to one of the realisations of the ordinary least squares method, e.g. Newton-Raphson algorithm, where however numerical integration of (3) will be replaced by approximation of the solution of the system (3) by an appropriately chosen element of some finite-dimensional function space. This seems to reduce considerably the computational load in comparison to the gradient-based optimisation methods which employ numerical integration of (3). However, it still appears to be computationally more intense than the two-step approach advocated in the present work.

The rest of the paper is organised as follows: in the next section we will detail the approach that we use and present its theoretical properties. In particular, we will

---

<sup>2</sup>The proofs of the main results in Brunel (2008) are incomplete and the main theorems seem to require further conditions in order to hold.

show that under appropriate conditions our two-step approach leads to a consistent estimator with a  $\sqrt{n}$  convergence rate, which is the best possible rate in regular parametric models<sup>3</sup>. Section 4 contains a discussion on the results obtained and possible extensions. The proofs of the main results are relegated to Section 5, while the Appendices contain some auxiliary statements.

### 3. RESULTS

First of all, we point out that in the present study we will be concerned with the asymptotic behaviour of an appropriate two-step estimator of  $\theta$  under a suitable sampling scheme. We will primarily be interested in intuitively understanding the behaviour of a relatively simple estimator of  $\theta$ , as well as in a clear presentation of the obtained results and the proofs. Consequently, the stated conditions will not always be minimal and can typically be relaxed at appropriate places.

We first define the sampling scheme.

**Condition 1.** *The observation times  $0 \leq t_1 < \dots < t_n \leq 1$  are deterministic and there exists a constant  $c_0 \geq 1$ , such that for all  $n$*

$$\max_{2 \leq i \leq n} |t_i - t_{i-1}| \leq \frac{c_0}{n}$$

*holds. Furthermore, there exists a constant  $c_1 > 0$ , such that for any interval  $A \subseteq [0, 1]$  of length  $|A|$  and all  $n \geq 1$  the inequality*

$$\frac{1}{n} \sum_{i=1}^n 1_{[t_i \in A]} \leq c_1 \max \left( |A|, \frac{1}{n} \right)$$

*holds.*

Hence, we observe the solution of the system (3) on the interval  $[0, 1]$ . Instead of  $[0, 1]$  we could have taken any other bounded interval. Conditions on  $t_1, \dots, t_n$  as in Condition 1 are typical in nonparametric regression, see e.g. Gasser and Müller (1984) and Section 1.7 in Tsybakov (2009), and they imply that  $t_1, \dots, t_n$  are distributed over  $[0, 1]$  in a sufficiently uniform manner. The most important example in which Condition 1 is satisfied, is when the observations are spaced equidistantly over  $[0, 1]$ , i.e. when  $t_j = j/n$  for  $j = 1, \dots, n$ . In this case one may take  $c_0 = 1$ . Notice that we do not necessarily assume that the initial condition  $x(0) = \xi$  is measured or is known. If it is, then it is incorporated into the observations and is used in the first step of the two-step estimation procedure.

**Condition 2.** *The random variables  $\epsilon_{ij}$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, d$ , from (4) are independent and are normally distributed with mean zero and finite variance  $\sigma_j^2$ .*

This assumption of Gaussianity of the  $\epsilon_{ij}$ 's may be dropped in various ways, as we will see below; see the note after Proposition 1 and Appendix B.

We next state a condition on the parameter set.

**Condition 3.** *The parameter set  $\Theta$  is a compact subset of  $\mathbb{R}^p$ .*

Compactness of  $\Theta$  allows one to put relatively weak conditions on the structure of the system (3), i.e. the function  $F$ .

---

<sup>3</sup>It is claimed in Liang and Wu (2008) that their two-step estimation procedure leads to a faster rate than  $\sqrt{n}$ , which is impossible. Indeed, their Theorem 2 and its proof are incorrect.



Just as the least squares method, see e.g. Jennrich (1969), the two-step approach also requires some regularity of the solutions of (3). In what follows, a derivative of any function  $f$  with respect to the variable  $y$  will be denoted by  $f'_y$ . For the second derivative of  $f$  with respect to  $y$  we will use the notation  $f''_{yy}$  with a similar convention for mixed derivatives.

**Condition 4.** *The following conditions hold:*

- (i) *the mapping  $F : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}^d$  from (3) is such that its second derivatives  $F''_{\theta\theta}, F''_{\theta x}, F''_{xx}$  are continuous;*
- (ii) *for all parameter values  $\theta \in \Theta$ , the solution  $x_\theta$  of (3) is defined on the interval  $[0, 1]$ ;*
- (iii) *for all parameter values  $\theta \in \Theta$ , the solution  $x_\theta$  of (3) is unique on  $[0, 1]$ ;*
- (iv) *for all parameter values  $\theta \in \Theta$ , the solution  $x_\theta$  of (3) is a  $C^\alpha$  function of  $t$  on the interval  $[0, 1]$  for some positive integer  $\alpha$ .*

Observe that Condition 4 (i) implies existence and uniqueness of the solution of (3) in some neighbourhood of 0. However, we want the existence and uniqueness to hold on the whole interval  $[0, 1]$  and therefore a priori require (ii) and (iii). Furthermore,  $\alpha \geq 2$  in (iv) is required when establishing appropriate asymptotic properties of nonparametric estimators of the solution  $x_\theta$  and its derivative, while  $\alpha \geq 3$  is needed in Propositions 3 and 4, and  $\alpha \geq 4$  in Theorem 1, respectively. Notice that for every  $\theta$  the solution  $x_\theta$  is of class  $C^\alpha$  in  $t$  in a neighbourhood of 0, provided for a given  $\theta$  the function  $F$  is of class  $C^\alpha$  in its first argument. However, we want this to hold on the whole interval  $[0, 1]$  and therefore require (iv). Since in the theory of chemical reaction networks, see for instance Sontag (2001), the components of  $F$  are usually polynomial or rational functions of  $x_1, \dots, x_d$  and  $\theta_1, \dots, \theta_p$ , the solution of (3) will be smooth enough in many examples and  $\alpha \geq 4$  is satisfied in a large number of practical examples. For the above-mentioned facts from the theory of ordinary differential equations see e.g. Chapter 2 in Arnold (1973). Also notice that the condition on  $F$  in Liang and Wu (2008), see Assumption C on p. 1573, puts severe restrictions on  $F$  and excludes e.g. quadratic nonlinearities of  $F$  in  $x_1, \dots, x_d$ . This, of course, has to be avoided.

Recall that our observations are  $Y_{ij} = x_{\theta_j}(t_i) + \epsilon_{ij}$  for  $i = 1, \dots, n, j = 1, \dots, d$ . We propose the following nonparametric estimator for  $x_{\theta_j}$ ,

$$(5) \quad \hat{x}_j(t) = \sum_{i=1}^n (t_i - t_{i-1}) \frac{1}{b} K\left(\frac{t - t_i}{b}\right) Y_{ij},$$

where  $K$  is a kernel function, while the number  $b = b_n > 0$  denotes a bandwidth that we take to depend on the sample size  $n$  in such a way that  $b_n \rightarrow 0$  as  $n \rightarrow \infty$ . In line with a traditional convention in kernel estimation theory, we will suppress the dependence of  $b_n$  on  $n$  in our notation, since no confusion will arise. When the  $t_i$ 's are equispaced, the estimator (5) can in essence be obtained by modifying the Nadaraya-Watson regression estimator, cf. p. 34 in Tsybakov (2009). It is usually called the Priestley-Chao estimator after the authors who first proposed it in Priestley and Chao (1972). As far as an estimator of  $x'_{\theta_j}(t)$  is concerned, we define it as the derivative of  $\hat{x}_j(t)$  with respect to  $t$ , choosing  $K$  as a differentiable function. Notice that the bandwidth  $b$  plays a role of regularisation parameter: too small a bandwidth results in an estimator with small bias, but large variance, while too large a bandwidth results in an estimator with small variance, but large bias,

see e.g. pp. 7–8 and 32 in Tsybakov (2009) for a relevant discussion. In principle one could use different bandwidth sequences for estimation of  $x_j$  for different  $j$ 's, but as can be seen from the proofs in Section 5, asymptotically this will not make a difference for an estimator of  $\theta$ . A similar remark applies to the use of different bandwidths for estimation of  $x_{\theta_j}$  and its derivative  $x'_{\theta_j}$ . Arguably, the estimator (5) is simple and there exist other estimators that may outperform it in certain respects in practice. However, as we will show later on, even such a simple estimator leads to a  $\sqrt{n}$ -consistent estimator of  $\theta$ .

Theoretical properties of the Priestley-Chao estimator were studied in Benedetti (1977), Priestley and Chao (1972), and Schuster and Yakowitz (1979). However, the first two papers do not cover its convergence in the  $L_\infty$  (supremum) norm, while the third one does not do it in the form required in the present work. Since this is needed in the sequel, we will supply the required statement, see Proposition 1 below.

To put things in a somewhat more general context than the one in our differential equations setting, consider the following regression model:

$$(6) \quad \begin{aligned} Y_i &= \mu(t_i) + \epsilon_i, \quad i = 1, \dots, n, \\ t_1, \dots, t_n &\text{ satisfy Condition 1,} \\ \epsilon_1, \dots, \epsilon_n &\text{ are i.i.d. Gaussian with } \mathbb{E}[\epsilon_i] = 0 \text{ and } \text{Var}[\epsilon_i] = \sigma^2 > 0. \end{aligned}$$

Our goal is to estimate the regression function  $\mu$  and its derivative  $\mu'$ . The estimator of  $\mu$  will be given by an expression similar to (5), namely

$$(7) \quad \hat{\mu}_n(t) = \sum_{i=1}^n (t_i - t_{i-1}) \frac{1}{b} K\left(\frac{t - t_i}{b}\right) Y_i,$$

while an estimator of  $\mu'$  will be given by  $\hat{\mu}'_n$ . We postulate the following condition on the kernel  $K$  for some strictly positive integer  $\alpha$ .

**Condition 5.** *The kernel  $K$  is symmetric and twice continuously differentiable, it has support within  $[-1, 1]$ , and it satisfies the integrability conditions:  $\int_{-1}^1 K(u) du = 1$  and  $\int_{-1}^1 u^\ell K(u) du = 0$  for  $\ell = 1, \dots, \alpha - 1$ .*

The following proposition holds.

**Proposition 1.** *Suppose the regression model (6) is given and Condition 5 holds. Fix  $0 < \delta < 1/2$ .*

(i) *If  $\mu$  is  $\alpha \geq 1$  times continuously differentiable and  $b \rightarrow 0$  as  $n \rightarrow \infty$ , then*

$$(8) \quad \sup_{t \in [\delta, 1-\delta]} |\hat{\mu}_n(t) - \mu(t)| = O_P \left( \sqrt{\left(b^\alpha + \frac{1}{nb^2}\right)^2 + \frac{\log n}{nb}} \right).$$

(ii) *If  $\mu$  is  $\alpha \geq 2$  times continuously differentiable and  $b \rightarrow 0$  as  $n \rightarrow \infty$ , then*

$$(9) \quad \sup_{t \in [\delta, 1-\delta]} |\hat{\mu}'_n(t) - \mu'(t)| = O_P \left( \sqrt{\left(b^{\alpha-1} + \frac{1}{nb^3}\right)^2 + \frac{\log n}{nb^3}} \right)$$

*is valid. In particular,  $\hat{\mu}_n$  and  $\hat{\mu}'_n$  are consistent on  $[\delta, 1 - \delta]$ , if  $nb^3 / \log n \rightarrow \infty$  holds additionally.*

Gaussianity of the  $\epsilon_i$ 's allows one to prove (8) and (9) by relatively elementary means. This assumption can be modified in various ways, for instance by assuming that the  $\epsilon_i$ 's are bounded, and we state and prove the corresponding modification of Proposition 1 in Appendix B, see Proposition 5. If we only assume that the  $\epsilon_i$ 's are i.i.d. with  $\mathbb{E}[\epsilon_i] = 0$  and  $\text{Var}[\epsilon_i] < \infty$ , then the analogues of (8) and (9) can be proved using the arguments from the proofs of Theorem 2 and Lemma 3 in Schuster and Yakowitz (1979). The rate of convergence will however be different and will lead to a stronger condition on  $\alpha$  in Theorem 1, which is the main result of the present paper. In general, normality of the measurement errors is a standard assumption in parameter estimation for systems of ordinary differential equations, see e.g. Girolami (2008), Hemker (1972), and Ramsay et al. (2007).

The following corollary is immediate from Proposition 1.

**Corollary 1.** *Under Conditions 1-5 we have for the estimator  $\hat{x}_j$*

$$(10) \quad \sup_{t \in [\delta, 1-\delta]} |\hat{x}_j(t) - x_{\theta_j}(t)| = O_P \left( \sqrt{\left(b^\alpha + \frac{1}{nb^2}\right)^2 + \frac{\log n}{nb}} \right)$$

and

$$(11) \quad \sup_{t \in [\delta, 1-\delta]} |\hat{x}'_j(t) - x'_{\theta_j}(t)| = O_P \left( \sqrt{\left(b^{\alpha-1} + \frac{1}{nb^3}\right)^2 + \frac{\log n}{nb^3}} \right),$$

provided  $\alpha \geq 2$  and  $b \rightarrow 0$  as  $n \rightarrow \infty$ . In particular,  $\hat{x}_j$  and  $\hat{x}'_j$  are consistent, if  $nb^3/\log n \rightarrow \infty$  holds additionally.

In the proof of Proposition 1 we need to use the continuous mapping theorem in order to prove convergence in probability of certain integrals of  $F$  and its derivatives with  $\hat{x}$  plugged in. This is where Corollary 1 is used.

Now that we have consistent (in an appropriate sense) estimators of  $x_{\theta_j}$  and  $x'_{\theta_j}$ , we can move to the second step in the construction of the two-step estimator of  $\theta$ . In particular, we define the estimator  $\hat{\theta}_n$  of  $\theta$  as

$$(12) \quad \begin{aligned} \hat{\theta}_n &= \operatorname{argmin}_{\eta \in \Theta} \int_0^1 \|\hat{x}'(t) - F(\hat{x}(t), \eta)\|^2 w(t) dt \\ &= \operatorname{argmin}_{\eta \in \Theta} M_{n,w}(\eta), \end{aligned}$$

where  $\|\cdot\|$  denotes the usual Euclidean norm and  $w$  is a weight function. We will refer to  $M_{n,w}(\eta)$  as a (random) criterion function. Since  $\Theta$  is compact and  $M_{n,w}$  is continuous in  $\eta$ , the minimiser  $\hat{\theta}_n$  always exists. The fact that  $\hat{\theta}_n$  is a measurable function of the observations  $Y_{ij}$  follows from Lemma 2 of Jennrich (1969). Notice that in Liang and Wu (2008) and Varah (1982) the criterion function is given by

$$\sum_{i=1}^n \|\tilde{x}'(t_i) - F(\tilde{x}(t_i), \eta)\|^2,$$

where  $\tilde{x}$  and  $\tilde{x}'$  are appropriate estimators of  $x_\theta$  and  $x'_\theta$ . However, in order to obtain a  $\sqrt{n}$ -consistent estimator of  $\theta$ , it is important to use an integral type criterion: the nonparametric estimators of  $x_\theta$  and  $x'_\theta$  have a slower convergence rate than  $\sqrt{n}$  and the latter has to be counterbalanced by some other means when estimating  $\theta$ . In light of this the choice of the weight function  $w$  also appears to be important. Furthermore, the observations  $Y_{ij}$  indirectly carry information on the entire curves

$x_{\theta_j}(t), t \in [0, 1]$ , and not only on the points  $x_{\theta_j}(t_i)$ . An integral type criterion allows one to exploit this fact in the second step of this two-step estimation procedure.

Introduce the asymptotic criterion

$$(13) \quad M_w(\eta) = \int_0^1 \|F(x_\theta(t), \theta) - F(x_\theta(t), \eta)\|^2 w(t) dt$$

corresponding to  $M_{n,w}$ . Observe that by Condition 4 it is bounded. Using Corollary 1 as a building block, one can show that the two-step estimator  $\hat{\theta}_n$  is consistent. To this end we will need the following condition on the weight function  $w$ .

**Condition 6.** *The weight function  $w$  is a nonnegative function that is continuously differentiable, is supported on the interval  $(\delta, 1 - \delta)$  for  $0 < \delta < 1/2$ , and is such that the Lebesgue measure of the set  $\{t : w(t) > 0\}$  is positive.*

The fact that  $w$  vanishes at the endpoints of the interval  $[\delta, 1 - \delta]$  and beyond, is needed to obtain a  $\sqrt{n}$ -consistent estimator of  $\theta$ . The condition that  $w$  is supported on  $(\delta, 1 - \delta)$  takes care of the boundary bias effects characteristic of the conventional kernel-type estimators, see e.g. Gasser and Müller (1984) for more information on this. Boundary effects in kernel estimation are usually remedied by using special boundary kernels, see e.g. van Es (1991), Gasser et al. (1985), Messer and Goldstein (1993). Using such a kernel, it can be expected that in our case as well the boundary effects will be eliminated and one may relax the requirement  $0 < \delta < 1/2$  from Condition 6 to allowing  $\delta = 0$ , i.e. to allowing  $w$  to be supported on  $(0, 1)$ . The condition that the weight function  $w$  is positive on a set with positive Lebesgue measure, is important for (14) to hold and in fact  $w(t) = 0$  a.e. would be a strange choice.

The following proposition is valid.

**Proposition 2.** *Suppose  $b \rightarrow 0$  and  $nb^3/\log n \rightarrow \infty$ . Under Conditions 1–6 and the additional identifiability condition*

$$(14) \quad \forall \varepsilon > 0, \inf_{\|\eta - \theta\| \geq \varepsilon} M_w(\eta) > M_w(\theta),$$

*we have  $\hat{\theta}_n \xrightarrow{P} \theta$ .*

The proposition is proved via a reasoning standard in the theory of M-estimation: we show that  $M_{n,w}$  converges to  $M_w$  and that the convergence is strong enough to imply the convergence of a minimiser  $\hat{\theta}_n$  of  $M_{n,w}$  to a minimiser  $\theta$  of  $M_w$ , cf. Section 5.2 of van der Vaart (1998). A necessary condition for (14) to hold is that  $x_\theta(\cdot) \neq x_{\theta'}(\cdot)$  for  $\theta \neq \theta'$ . The latter is a minimal assumption for statistical identifiability of parameter  $\theta$ . The identifiability condition (14) is common in the theory of M-estimation, see Theorem 5.7 of van der Vaart (1998). It means that  $\theta$  is a point of minimum of  $M_w(\eta)$  and that it is a *well-separated* point of minimum. The most trivial example with this condition satisfied is when  $d = p = 1$  and  $x'(t) = \theta x(t)$  hold with initial condition  $x(0) = \xi$ , where  $\xi \neq 0$ . Observe that since  $\Theta$  is compact and  $M_w$  is continuous, uniqueness of a minimiser of  $M_w$  implies (14), cf. Exercise 27 on p. 84 of van der Vaart (1998).

In practice (14) might be difficult to check globally and one might prefer to concentrate on a simpler local condition: if the first order condition  $[dM_w(\eta)/d\eta]_{\eta=\theta} = 0$  holds and if the Hessian matrix  $H(\eta) = (\partial^2 M_w(\eta)/\partial \eta_i \partial \eta_j)_{i,j}$  of  $M_w$  is strictly

positive definite at  $\theta$ , then (14) will be satisfied for  $\eta \in \Theta$  restricted to some neighbourhood of  $\theta$ , because  $M_w$  will have a local minimum at such  $\theta$  and a neighbourhood around it can be taken to be compact with small enough diameter, so that (14) holds for  $\eta$  restricted to this neighbourhood. The conclusion of the theorem will then hold for the parameter set restricted to this neighbourhood of  $\theta$ .

In a statement analogous to Proposition 2, Brunel (2008) requires that the solutions of (3) belong to a compact set  $\mathcal{K}$  for all  $\theta$  and  $t$  and that  $F$  from (1) is Lipschitz in its first argument  $x$  for  $x$  restricted to this compact  $\mathcal{K}$  uniformly in  $\theta \in \Theta$ . It is also assumed that the nonparametric estimators  $\hat{x}_n(t)$  belong a.s. to  $\mathcal{K}$  for all  $n$  and  $t$ . However, the latter typically will not hold for linear smoothers, see Definition 1.7 in Tsybakov (2009), which constitute the most popular choice of nonparametric regression estimators in practice. For instance, local polynomial estimators, see Section 1.6 in Tsybakov (2009), projection estimators, see Section 1.7 in Tsybakov (2009), or the Gasser-Müller estimator, see Gasser and Müller (1984), are all examples of linear smoothers. Hence we prefer to avoid this condition altogether, although this somewhat complicates the proof.

Under the conditions in this section it turns out that the estimator  $\hat{\theta}_n$  is not merely a consistent estimator, but a  $\sqrt{n}$ -consistent estimator of  $\theta$ , in the sense of (18) below. This result follows in essence from the fact that up to a higher order term the difference  $\hat{\theta}_n - \theta$  can be represented as the difference of the images of  $\hat{x}$  and  $x_\theta$  under a certain linear mapping, cf. (30). It is known that even though nonparametric curve estimators cannot usually attain the  $\sqrt{n}$  convergence rate, see e.g. Chapters 1 and 2 of Tsybakov (2009), extra smoothness often coming from the structure of linear functionals allows one to construct in many cases  $\sqrt{n}$ -consistent estimators of these functionals via plugging in nonparametric estimators, see e.g. Bickel and Ritov (2003) and Goldstein and Messer (1992) for more information. The variance of such plug-in estimators can often be proven to be of order  $n^{-1}$ , while the squared bias can be made of order  $n^{-1}$  by undersmoothing, i.e. selecting the smoothing parameter smaller than what is an optimal choice in nonparametric curve estimation when the object of interest is a curve itself, cf. Goldstein and Messer (1992). Precisely this happens in our case as well: if the mean integrated squared error is used as a performance criterion of a nonparametric estimator, then under our conditions the optimal bandwidth for estimation of  $x_\theta$  is of order  $n^{-1/(2\alpha+1)}$ , whereas the optimal bandwidth for estimation of  $\theta$  is in fact smaller, see Theorem 1 below. Note that this is a different approach than the one in Bickel and Ritov (2003), where it is assumed that nonparametric estimators attain the minimax rate of convergence and the  $\sqrt{n}$ -rate for estimation of a functional in concrete examples, if possible, is achieved by different means exploiting extra smoothness coming from the structure of a functional, see e.g. the first example in Section 2 there. In many cases it can be proved that such plug-in type estimators are efficient, see Bickel and Ritov (2003). Notice, however, that in our case this will not imply that  $\hat{\theta}_n$  is efficient.

First we will provide an asymptotic representation for the difference  $\hat{\theta}_n - \theta$ .

**Proposition 3.** *Let  $\theta$  be an interior point of  $\Theta$ . Suppose that the conditions of Proposition 2 hold and let the matrix  $J_\theta$  defined by*

$$(15) \quad J_\theta = \int_{\delta}^{1-\delta} (F'_\theta(x_\theta(t), \theta))^T F'_\theta(x_\theta(t), \theta) w(t) dt$$

be nonsingular. Fix  $\alpha \geq 3$ . If  $b \asymp n^{-\gamma}$  holds for  $1/(4\alpha - 4) < \gamma < 1/6$ , then

$$(16) \quad \hat{\theta}_n - \theta = O_P(J_\theta^{-1}(\Gamma(\hat{x}) - \Gamma(x_\theta))) + o_P(n^{-1/2})$$

is valid with the mapping  $\Gamma$  given by

$$(17) \quad \Gamma(z) = \int_\delta^{1-\delta} \left\{ (F'_\theta(x_\theta(t), \theta))^T F'_x(x_\theta(t), \theta) w(t) - \frac{d}{dt} [(F'_\theta(x_\theta(t), \theta))^T w(t)] \right\} z(t) dt.$$

With the above result in mind, in order to complete the study of the asymptotics of  $\hat{\theta}_n$ , it remains to study the mapping  $\Gamma$ . Clearly, it suffices to study the asymptotic behaviour of

$$\Delta(\hat{\mu}_n) - \Delta(\mu) = \int_{\mathbb{R}} v(t)k(t)\hat{\mu}_n(t)dt - \int_{\mathbb{R}} v(t)k(t)\mu(t)dt,$$

where  $v$  is a known function that satisfies appropriate assumptions, while  $k$  stands either for  $w$  or its derivative  $w'$ . The next proposition deals with the asymptotics of  $\Delta(\hat{\mu}_n) - \Delta(\mu)$ .

**Proposition 4.** *Under Conditions 5 and 6 and for any continuous function  $v$  it holds in the regression model (6) that*

$$\Delta(\hat{\mu}_n) - \Delta(\mu) = O_P(n^{-1/2}),$$

provided  $\mu$  is  $\alpha \geq 3$  times differentiable and the bandwidth  $b$  is chosen such that  $b \asymp n^{-\gamma}$  holds for  $1/(2\alpha) \leq \gamma \leq 1/4$ .

Our main result is a simple consequence of Propositions 3 and 4.

**Theorem 1.** *Let  $\theta$  be an interior point of  $\Theta$ . Assume that Conditions 1–6 together with (14) hold and that (15) is nonsingular. Fix  $\alpha \geq 4$ . If the bandwidth  $b$  is such that  $b \asymp n^{-\gamma}$  holds for  $1/(2\alpha) < \gamma < 1/6$ , then*

$$(18) \quad \sqrt{n}(\hat{\theta}_n - \theta) = O_P(1)$$

is valid.

Thus any bandwidth sequences satisfying the conditions in Theorem 1 are optimal, in the sense that they lead to estimators with similar asymptotic behaviour. In particular, each of such bandwidth sequences ensures a  $\sqrt{n}$  convergence rate of  $\hat{\theta}_n$ . Consequently, dependence of the asymptotic properties of the estimator  $\hat{\theta}_n$  on the bandwidth is less critical than it typically is in nonparametric curve estimation. Notice that the condition  $\alpha \geq 4$  in Theorem 1 is needed in order to make the conditions in Propositions 3 and 4 compatible.

#### 4. DISCUSSION

The main result of the paper, Theorem 1, is that under certain conditions for systems of ordinary differential equations parameter estimation at the  $\sqrt{n}$  rate is possible *without* employing numerical integration. Although we have shown this in the case when in the first step of the two-step procedure a particular kernel-type estimator is used, it may be expected that a similar result holds for other nonparametric estimators. In practice for small or moderate sample sizes it might be advantageous to use more sophisticated nonparametric estimators than the Priestley-Chao estimator, but asymptotically this does not make a difference.

Once a  $\sqrt{n}$ -consistent estimator  $\hat{\theta}_n$  of  $\theta$  is available, one might ask for more, namely if one can construct an estimator that is asymptotically equivalent to the ordinary least squares estimator (2) or that is semiparametrically efficient. It is expected that this can be achieved without repeated numerical integration of (1) by using  $\hat{\theta}_n$  as a starting point and performing a one-step Newton-Raphson type procedure; see e.g. Section 7.8 of Bickel et al. (1998) or Chapter 25 of van der Vaart (1998). We intend to address this issue of efficient and ordinary least squares estimation in a separate publication.

Doubtless, the main challenge in implementing the two-step estimation procedure lies in selecting the smoothing parameter  $b$ . This is true for any two-step procedure, e.g. the one based on the regression splines as in Brunel (2008) or the local polynomial estimator as in Liang and Wu (2008), and not only for our specific estimator. Observations that we supply below apply in principle to any two-step estimator and not only to the specific one considered in the present work. Hence they are of general interest.

Some attention has been paid in the literature to the selection of the smoothing parameter in the context of parameter estimation for ordinary differential equations. The considered options range from subjective choices and smoothing by hand to more advanced possibilities. Perhaps the simplest solution would be to assume that the targets of the estimation procedure are  $x_{\theta_j}$ ,  $j = 1, \dots, d$ , and to select  $b$  (a different one for every component  $x_{\theta_j}$ ) via a cross-validation procedure, see e.g. Section 5.3 in Wasserman (2006) for a description of cross-validation techniques in the context of nonparametric regression. This should produce reasonable results, at least for relatively large sample sizes, cf. simulation examples considered in Brunel (2008). However, it is clear from Theorem 1 that despite its simplicity, such a choice of  $b$  will be suboptimal. One other possibility for practical bandwidth selection is nothing else but a variation on the plug-in bandwidth selection method as described e.g. in Jones et al. (1996): if one computes the mean squared error of  $\hat{\theta}_n$ , one can see from the proof in Section 5 that the terms that depend on the bandwidth  $b$  are lower order terms in the expansion of the mean squared error. One can then minimise with respect to  $b$  a bound on these lower order terms. A minimiser, say  $b^*$ , will depend on the unknown true parameter  $\theta$ , also via  $x_\theta$  and  $x'_\theta$ , as well as on the error variances  $\sigma_1^2, \dots, \sigma_d^2$ . However,  $\theta$ ,  $x_\theta$ , and  $x'_\theta$  can be re-estimated via  $\hat{\theta}_n$ ,  $\hat{x}$ , and  $\hat{x}'$  using a different, pilot bandwidth  $\hat{b}$ . Of course, instead of  $\hat{x}$  and  $\hat{x}'$  the use of any other nonparametric estimators of a regression function and its derivative, e.g. local polynomial estimators, see Section 1.6 of Tsybakov (2009), or the Gasser-Müller estimator, see Gasser and Müller (1984), is also a valid option. Error term variances can be estimated via one of the methods described in Hall and Marron (1990) or Section 5.6 of Wasserman (2006). Once the pilot estimators of  $\theta$ ,  $x_\theta$ , and  $x'_\theta$  together with estimators of  $\sigma_1^2, \dots, \sigma_d^2$  are available, these can be plugged back into  $b^*$  and in this way one obtains a bandwidth  $\hat{b}$  that estimates the optimal bandwidth  $b^*$ . The final step would be computation of  $\hat{\theta}_n$  with a new bandwidth  $\hat{b}$ . Unfortunately, this method leads to extremely cumbersome expressions and furthermore, since we are minimising an upper bound on numerous remainder terms, it will probably tend to oversmooth, i.e. produce a bandwidth  $b$  larger than required. Moreover, the plug-in approach in general is subject to some controversy having both supporters and critics, see e.g. Loader (1999) and references therein. An alternative to the plug-in approach might be an approach based on one of the resampling methods:

cross-validation, jackknife, or bootstrap. Theoretical analysis of the properties of such bandwidth selectors is a rather nontrivial task. Also a thorough simulation study is needed before the practical value of different bandwidth selection methods can be assessed. We do not address these issues here.

The next observation of this section concerns the numerical computation of the two-step estimator. The kernel-type nonparametric regression estimates of  $x_{\theta_j}$ ,  $j = 1, \dots, d$ , can be quickly evaluated on any regular grid of points  $0 \leq s_1 \leq \dots \leq s_m$ , via techniques using the Fast Fourier Transform (FFT) similar to those described in Appendix D of Wand and Jones (1995). Furthermore, in the second step of the two-step estimation procedure the criterion function  $M_{n,w}$  can be approximated by a finite sum by discretising the integral in its definition. If  $F$  is linear in  $\theta_1, \dots, \theta_p$ , then as already observed in Varah (1982), see pp. 29 and 31, cf. p. 1262 in Brunel (2008) and p. 1573 in Liang and Wu (2008), this will lead to a weighted linear least squares problem, which can be solved in a routine fashion without using e.g. random search methods. This is a great simplification in comparison to the ordinary least squares estimator, which moreover will still tend to get trapped in local minima of the least squares criterion function despite the fact that  $F$  is linear in its parameters.

We conclude this section by mentioning one way for possible extension of the two-step method described in the present work. The last two decades have seen much interest and research in estimation methods for models with high-dimensional parameter spaces under assumptions of sparsity. Roughly speaking this means that even though the parameter indexing the model is a high- or possibly infinite-dimensional vector, many of its components are known to be zero. Within the context of systems of differential equations this appears to be the case for the so-called S-systems, see Voit (2000) for a comprehensive treatment of the theory of S-systems. The number of parameters for a  $d$ -dimensional S-system is given by  $2d(d+1)$  and it thus grows fairly fast with the dimension  $d$ . On the other hand practical experience indicates that many of the parameters in S-systems are equal to zero. One possible way of incorporating this information into the two-step estimation procedure would be to use the modified criterion function

$$\widetilde{M}_w(\eta) = \int_0^1 \|\hat{x}'(t) - F(\hat{x}(t), \eta)\|^2 w(t) dt + \lambda \text{pen}[\eta],$$

where  $\text{pen}[\eta]$  is a penalty for the size of  $\eta$  when choosing a specific parameter value  $\eta$ , while  $\lambda > 0$  is a tuning parameter that quantifies the degree of penalisation. A typical choice for  $\text{pen}[\eta]$  would be the  $L_1$  norm of  $\eta$ ,

$$\text{pen}[\eta] = \sum_{j=1}^p |\eta_j|.$$

A similar idea, but within the nonlinear least squares framework relying on numerical integration of (3), is explored in Kikuchi et al. (2003) on a model with simulated data. The method proposed there, is called ‘pruning’ by the authors. The authors however do not perform a study of the asymptotics of their estimator nor do they propose a practical method for selection of the penalty parameter.

We intend to perform a more practically oriented study exploring these ideas in a separate publication.



## 5. PROOFS

We will use the symbol  $\lesssim$ , meaning less or equal up to a universal constant independent of index  $n$ . The symbol  $\asymp$  will denote the fact that two sequences of real numbers are asymptotically of the same order.

*Proof of Proposition 1.* We first prove (8). For any positive  $\varepsilon$  by Chebyshev's inequality we have

$$(19) \quad \begin{aligned} P \left( \sup_{t \in [\delta, 1-\delta]} |\hat{\mu}_n(t) - \mu(t)| > \varepsilon \right) &\leq \frac{2}{\varepsilon^2} \left\{ \sup_{t \in [\delta, 1-\delta]} |\mathbb{E}[\hat{\mu}_n(t)] - \mu(t)|^2 \right. \\ &\quad \left. + \mathbb{E} \left[ \sup_{t \in [\delta, 1-\delta]} |\hat{\mu}_n(t) - \mathbb{E}[\hat{\mu}_n(t)]|^2 \right] \right\} \\ &= \frac{2}{\varepsilon^2} (T_1 + T_2). \end{aligned}$$

By (36) we can write

$$\mathbb{E}[\hat{\mu}_n(t)] - \mu(t) = \int_0^1 \mu(s) \frac{1}{b} K \left( \frac{t-s}{b} \right) ds - \mu(t) + O \left( \frac{1}{nb^2} \right).$$

For all  $n$  large enough, we have  $b \leq \delta$ , because  $b \rightarrow 0$ . Then for all such  $n$ , if  $t \in [\delta, 1-\delta]$ , a standard argument (cf. p. 6 in Tsybakov (2009)), namely Taylor's formula up to order  $\alpha$  applied to  $\mu$  and the moment conditions on the kernel  $K$  formulated in Condition 5, yields

$$(20) \quad \sup_{t \in [\delta, 1-\delta]} |\mathbb{E}[\hat{\mu}_n(t)] - \mu(t)| \leq b^\alpha \frac{\|\mu^{(\alpha)}\|_\infty}{\alpha!} \int_{-1}^1 |u^\alpha K(u)| du + O \left( \frac{1}{nb^2} \right).$$

Next we turn to  $T_2$ . With argumentation similar to that in the proof of Theorem 1.8 of Tsybakov (2009) and setting  $S_i(t) = (t_i - t_{i-1})/b K((t - t_i)/b)$ ,  $N = n^2$ , and  $s_j = j/N$ , for  $j = 1, \dots, N$ , we have

$$\begin{aligned} A &= \sup_{t \in [\delta, 1-\delta]} |\hat{\mu}_n(t) - \mathbb{E}[\hat{\mu}_n(t)]| \\ &= \sup_{t \in [\delta, 1-\delta]} \left| \sum_{i=1}^n \epsilon_i S_i(t) \right| \\ &\leq \max_{1 \leq j \leq N} \left| \sum_{i=1}^n \epsilon_i S_i(s_j) \right| + \sup_{t, t': |t-t'| \leq N^{-1}} \left| \sum_{i=1}^n \epsilon_i (S_i(t) - S_i(t')) \right|. \end{aligned}$$

By the mean value theorem and Condition 1 the inequality

$$|S_i(t) - S_i(t')| \lesssim \|K'\|_\infty \frac{1}{nb^2} |t - t'|$$

holds for any  $t, t' \in \mathbb{R}$ , where  $\|K'\|_\infty$  is finite. Hence by the  $c_2$ -inequality

$$\begin{aligned} A^2 &\leq \left( \max_{1 \leq j \leq N} \left| \sum_{i=1}^n \epsilon_i S_i(s_j) \right| + \sup_{t, t': |t-t'| \leq N^{-1}} \left| \sum_{i=1}^n \epsilon_i (S_i(t) - S_i(t')) \right| \right)^2 \\ &\lesssim \max_{1 \leq j \leq N} |Z_j|^2 + \frac{\|K'\|_\infty^2}{n^2 b^4 N^2} \left( \sum_{i=1}^n |\epsilon_i| \right)^2, \end{aligned}$$

where  $Z_j = \sum_{i=1}^n \epsilon_i S_i(s_j)$ . Therefore we have

$$\mathbb{E}[A^2] \leq \mathbb{E} \left[ \max_{1 \leq j \leq N} |Z_j|^2 \right] + \frac{\|K'\|_\infty^2}{n^2 b^4 N^2} \mathbb{E} \left[ \left( \sum_{i=1}^n |\epsilon_i| \right)^2 \right].$$

Notice that

$$\frac{1}{n^2 b^4 N^2} \mathbb{E} \left[ \left( \sum_{i=1}^n |\epsilon_i| \right)^2 \right] \leq \frac{\mathbb{E}[\epsilon_1^2]}{N^2 b^4} = \frac{\sigma^2}{n^4 b^4} = o\left(\frac{1}{nb}\right).$$

Moreover, we have

$$\begin{aligned} \mathbb{E}[Z_j^2] &= \sum_{i=1}^n \sigma^2 (t_i - t_{i-1})^2 \left( \frac{1}{b} K \left( \frac{t_i - s_j}{b} \right) \right)^2 \\ &\lesssim \frac{\sigma^2 \|K\|_\infty^2}{n^2 b^2} \sum_{i=1}^n \mathbb{1}_{|t_i - s_j| \leq b} \\ &\leq \frac{1}{nb} c_1 \sigma^2 \|K\|_\infty^2 \max\left(2, \frac{1}{nb}\right), \end{aligned}$$

where the last inequality follows from Condition 1. Since the  $Z_j$ 's, being a linear combination of independent Gaussian random variables, are themselves Gaussian, Corollary 1.3 of Tsybakov (2009) and the fact that  $N = n^2$  then entail

$$\mathbb{E} \left[ \max_{1 \leq j \leq N} |Z_j|^2 \right] = O\left(\frac{\log N}{nb}\right) = O\left(\frac{\log n}{nb}\right).$$

Hence

$$(21) \quad \mathbb{E}[A^2] = O((\log n)/(nb)).$$

Taking

$$\varepsilon = M \sqrt{\left(b^\alpha + \frac{1}{nb^2}\right)^2 + \frac{\log n}{nb}}$$

with an appropriate constant  $M$  yields (8) by (19), (20), and (21).

As far as the proof of (9) is concerned, it is very much similar to the proof of (8) and is therefore omitted. This completes the proof of the proposition.  $\square$

*Proof of Proposition 2.* From the definition of  $M_{n,w}(\eta)$  and  $M_w(\eta)$ , the elementary inequality

$$\left| \|a_1\|^2 - \|a_2\|^2 \right| \leq \|a_1 - a_2\| (\|a_1\| + \|a_2\|)$$

and the Cauchy-Schwarz inequality we have

$$\begin{aligned} & |M_{n,w}(\eta) - M_w(\eta)| \\ & \leq \left\{ \int_\delta^{1-\delta} \|\hat{x}'(t) - F(x_\theta(t), \theta) + F(x_\theta(t), \eta) - F(\hat{x}(t), \eta)\|^2 w(t) dt \right\}^{1/2} \\ (22) \quad & \times \left\{ \int_\delta^{1-\delta} (\|\hat{x}'(t) - F(\hat{x}(t), \eta)\| + \|F(x_\theta(t), \theta) - F(x_\theta(t), \eta)\|)^2 w(t) dt \right\}^{1/2} \\ & = \sqrt{T_1} \sqrt{T_2}. \end{aligned}$$

For  $T_1$  we have that

$$(23) \quad T_1 \leq 2 \int_{\delta}^{1-\delta} \|\hat{x}'(t) - F(x_{\theta}(t), \theta)\|^2 w(t) dt \\ + 2 \int_{\delta}^{1-\delta} \|F(x_{\theta}(t), \eta) - F(\hat{x}(t), \eta)\|^2 w(t) dt.$$

By (11) it holds that

$$(24) \quad \sup_{\eta \in \Theta} \int_{\delta}^{1-\delta} \|\hat{x}'(t) - F(x_{\theta}(t), \theta)\|^2 w(t) dt \\ = \int_{\delta}^{1-\delta} \|\hat{x}'(t) - x'_{\theta}(t)\|^2 w(t) dt \\ \leq \sum_{i=1}^d \sup_{t \in [\delta, 1-\delta]} |\hat{x}'_i(t) - x'_{i,\theta}(t)|^2 \int_{\delta}^{1-\delta} w(t) dt \\ \xrightarrow{P} 0.$$

Moreover, by Lemma 3 from Appendix A we obtain that

$$(25) \quad \sup_{\eta \in \Theta} \int_{\delta}^{1-\delta} \|F(\hat{x}(t), \eta) - F(x_{\theta}(t), \eta)\|^2 w(t) dt \xrightarrow{P} 0.$$

Furthermore,  $T_2 = O_P(1)$  as  $n \rightarrow \infty$ , because

$$(26) \quad \sup_{\eta \in \Theta} \int_{\delta}^{1-\delta} \|F(x_{\theta}(t), \theta) - F(x_{\theta}(t), \eta)\|^2 w(t) dt < \infty$$

by compactness of  $\Theta$  and Condition 4, and

$$(27) \quad \sup_{\eta \in \Theta} \int_{\delta}^{1-\delta} \|\hat{x}'(t) - F(\hat{x}(t), \eta)\|^2 w(t) dt \xrightarrow{P} 0$$

by Lemma 4 from Appendix A. Combination of (22)–(27) implies that

$$\sup_{\eta \in \Theta} |M_{n,w}(\eta) - M_w(\eta)| \xrightarrow{P} 0.$$

The statement of the proposition then follows from this fact, the identifiability condition (14), and Theorem 5.7 of van der Vaart (1998).  $\square$

*Proof of Proposition 3.* We interpret the derivative of a one-dimensional function of  $\theta$  as a row  $p$ -vector of partial derivatives and we denote the  $d \times p$ -matrix of partial derivatives  $\frac{\partial}{\partial \theta_j} F_i(x, \theta)$ ,  $i = 1, \dots, d$ ,  $j = 1, \dots, p$ , by  $F'_{\theta}(x, \theta)$ .

We have

$$\frac{d}{d\theta} \|\hat{x}'(t) - F(\hat{x}(t), \theta)\|^2 = -2(\hat{x}'(t) - F(\hat{x}(t), \theta))^T F'_{\theta}(\hat{x}(t), \theta).$$

With this in mind and interchanging the order of integration and differentiation, we find that the derivative of  $M_{n,w}$  from (12) with respect to  $\theta$  is given by

$$-2 \int_{\delta}^{1-\delta} (\hat{x}'(t) - F(\hat{x}(t), \theta))^T F'_{\theta}(\hat{x}(t), \theta) w(t) dt.$$

Since  $\theta$  is an interior point of  $\Theta$ , there exists  $\varepsilon > 0$ , such that  $\text{ball}(\theta, \varepsilon)$ , the open ball of radius  $\varepsilon$  around  $\theta$ , is contained in  $\Theta$ . Take

$$G_n = \{|\hat{\theta}_n - \theta| < \varepsilon/2\}$$

and notice that by consistency of  $\hat{\theta}_n$  we have  $P(G_n) \rightarrow 1$  as  $n \rightarrow \infty$ . If  $\hat{\theta}_n$  is a point of minimum of  $M_{n,w}$ , then necessarily

$$1_{G_n} \int_{\delta}^{1-\delta} (\hat{x}'(t) - F(\hat{x}(t), \hat{\theta}_n))^T F'_{\theta}(\hat{x}(t), \hat{\theta}_n) w(t) dt = 0,$$

where 0 at the righthand side denotes now a row  $p$ -vector with all its entries equal to zero. The latter display can be rearranged as

$$\begin{aligned} 1_{G_n} \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T \times \{(\hat{x}'(t) - x'_{\theta}(t)) \\ + (F(x_{\theta}(t), \theta) - F(\hat{x}(t), \theta)) + (F(\hat{x}(t), \theta) - F(\hat{x}(t), \hat{\theta}_n))\} w(t) dt = 0, \end{aligned}$$

where now 0 on the righthand side denotes a column  $p$ -vector with its entries equal to zero. Note that we have

$$F(\hat{x}(t), \theta) - F(\hat{x}(t), \hat{\theta}_n) = \int_0^1 F'_{\theta}(\hat{x}(t), \hat{\theta}_n + \lambda(\theta - \hat{\theta}_n)) d\lambda (\theta - \hat{\theta}_n).$$

Hence

$$\begin{aligned} (28) \quad & 1_{G_n} \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T \int_0^1 F'_{\theta}(\hat{x}(t), \hat{\theta}_n + \lambda(\theta - \hat{\theta}_n)) d\lambda w(t) dt (\hat{\theta}_n - \theta) \\ & = 1_{G_n} \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T (\hat{x}'(t) - x'_{\theta}(t)) w(t) dt \\ & \quad + 1_{G_n} \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T (F(x_{\theta}(t), \theta) - F(\hat{x}(t), \theta)) w(t) dt \end{aligned}$$

holds. By the fact that  $\hat{x}$  converges in probability as a random element on  $[\delta, 1 - \delta]$  to  $x_{\theta}$ , see (10), consistency of  $\hat{\theta}_n$ , continuity of  $F'_{\theta}$ , continuity of integration and the continuous mapping theorem, see Theorem 18.11 in van der Vaart (1998), we have

$$\begin{aligned} (29) \quad & \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T \int_0^1 F'_{\theta}(\hat{x}(t), \hat{\theta}_n + \lambda(\theta - \hat{\theta}_n)) d\lambda w(t) dt \\ & \xrightarrow{P} \int_{\delta}^{1-\delta} (F'_{\theta}(x_{\theta}(t), \theta))^T F'_{\theta}(x_{\theta}(t), \theta) w(t) dt = J_{\theta}, \end{aligned}$$

where  $J_{\theta}$  is nonsingular by assumption (15). Therefore, the asymptotic behaviour of  $\hat{\theta}_n - \theta$  is given by

$$\begin{aligned} (30) \quad & J_{\theta}^{-1} \left( \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T (\hat{x}'(t) - x'_{\theta}(t)) w(t) dt \right. \\ & \quad \left. + \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T (F(x_{\theta}(t), \theta) - F(\hat{x}(t), \theta)) w(t) dt \right). \end{aligned}$$

It thus remains to be shown that this expression in fact reduces to the righthand side of (16). First of all, notice that

$$\begin{aligned}
& \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T (\hat{x}'(t) - x'_{\theta}(t)) w(t) dt \\
&= \int_{\delta}^{1-\delta} (F'_{\theta}(x_{\theta}(t), \theta))^T (\hat{x}'(t) - x'_{\theta}(t)) w(t) dt \\
(31) \quad &+ \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta))^T (\hat{x}'(t) - x'_{\theta}(t)) w(t) dt \\
&= - \int_{\delta}^{1-\delta} \left( \frac{d}{dt} [F'_{\theta}(x_{\theta}(t), \theta) w(t)] \right)^T (\hat{x}(t) - x_{\theta}(t)) dt \\
&+ \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta))^T (\hat{x}'(t) - x'_{\theta}(t)) w(t) dt,
\end{aligned}$$

where the last equality follows by integration by parts and the fact that  $w(\delta) = w(1-\delta) = 0$ . The first term at the righthand side of (31) appears also in the leading term  $\Gamma(\hat{x}) - \Gamma(x_{\theta})$  of (16). We will now show that the other term at the righthand side of (31) is negligible, i.e.

$$\int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta))^T (\hat{x}'(t) - x'_{\theta}(t)) w(t) dt = o_P(n^{-1/2}).$$

By the Cauchy-Schwarz inequality

$$\begin{aligned}
& \left\| \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta))^T (\hat{x}'(t) - x'_{\theta}(t)) w(t) dt \right\| \\
&\leq \left\{ \int_{\delta}^{1-\delta} \|F'_{\theta}(\hat{x}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta)\|^2 w(t) dt \right\}^{1/2} \\
&\quad \times \left\{ \int_{\delta}^{1-\delta} \|\hat{x}'(t) - x'_{\theta}(t)\|^2 w(t) dt \right\}^{1/2},
\end{aligned}$$

where  $\|\cdot\|$  denotes the Frobenius or the Hilbert-Schmidt norm of a matrix (recall that it is submultiplicative). By (11) we have

$$\left\{ \int_{\delta}^{1-\delta} \|\hat{x}'(t) - x'_{\theta}(t)\|^2 w(t) dt \right\}^{1/2} = O_P(1) \sqrt{\left( b^{\alpha-1} + \frac{1}{nb^3} \right)^2 + \frac{\log n}{nb^3}}.$$

Furthermore,

$$\begin{aligned}
& \int_{\delta}^{1-\delta} \|F'_{\theta}(\hat{x}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta)\|^2 w(t) dt \\
(32) \quad &\leq 2 \int_{\delta}^{1-\delta} \|F'_{\theta}(\hat{x}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \hat{\theta}_n)\|^2 w(t) dt \\
&+ 2 \int_{\delta}^{1-\delta} \|F'_{\theta}(x_{\theta}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta)\|^2 w(t) dt \\
&= 2T_1 + 2T_2.
\end{aligned}$$

Denote  $F'_{\theta}(x, \theta) = A(x, \theta) = (a_{i,j}(x, \theta))_{i,j}$ . For  $T_1$  we have

$$\begin{aligned}
T_1 &= \sum_{i,j} \int_{\delta}^{1-\delta} (a_{i,j}(\hat{x}(t), \hat{\theta}_n) - a_{i,j}(x_{\theta}(t), \hat{\theta}_n))^2 w(t) dt \\
&= \sum_{i,j} \int_{\delta}^{1-\delta} \left( \int_0^1 \frac{\partial}{\partial x} a_{i,j}(x_{\theta}(t) + \lambda(\hat{x}(t) - x_{\theta}(t)), \hat{\theta}_n) d\lambda (\hat{x}(t) - x_{\theta}(t)) \right)^2 w(t) dt \\
&\leq \left( \sup_{t \in [\delta, 1-\delta]} \|\hat{x}(t) - x_{\theta}(t)\|^2 \right) \\
&\quad \times \sum_{i,j} \int_{\delta}^{1-\delta} \int_0^1 \left\| \frac{\partial}{\partial x} a_{i,j}(x_{\theta}(t) + \lambda(\hat{x}(t) - x_{\theta}(t)), \hat{\theta}_n) \right\|^2 d\lambda w(t) dt.
\end{aligned}$$

By (10), as well as consistency of  $\hat{\theta}_n$ , Condition 4 and the continuous mapping theorem, the righthand side in the last inequality is of order

$$O_P(1) \left\{ \left( b^{\alpha} + \frac{1}{nb^2} \right)^2 + \frac{\log n}{nb} \right\}.$$

By a similar argument, the inequality

$$\begin{aligned}
T_2 &= \int_{\delta}^{1-\delta} \|F'_{\theta}(x_{\theta}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta)\|^2 w(t) dt \\
&\leq \|\hat{\theta}_n - \theta\|^2 \sum_{i,j} \int_{\delta}^{1-\delta} \int_0^1 \left\| \frac{\partial}{\partial \theta} a_{i,j}(x_{\theta}(t), \theta + \lambda(\hat{\theta}_n - \theta)) \right\|^2 d\lambda w(t) dt
\end{aligned}$$

holds. Here with some natural abuse of notation we first differentiate  $a_{i,j}$  with respect to its second argument  $\theta$  and only afterwards evaluate the obtained derivative at  $x_{\theta}(t)$  and  $\theta + \lambda(\hat{\theta}_n - \theta)$ . Since the integrals in the last inequality in the above display are bounded in probability, we then get

$$(33) \quad \left\{ \int_{\delta}^{1-\delta} \|F'_{\theta}(x_{\theta}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta)\|^2 w(t) dt \right\}^{1/2} = O_P(\|\hat{\theta}_n - \theta\|).$$

Now notice that (30) yields

$$\begin{aligned}
\|\hat{\theta}_n - \theta\| &\leq O_P(1) \left( \left\| \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T (\hat{x}'(t) - x'_{\theta}(t)) w(t) dt \right\| \right. \\
&\quad \left. + \left\| \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T (F(x_{\theta}(t), \theta) - F(\hat{x}(t), \theta)) w(t) dt \right\| \right).
\end{aligned}$$

The Cauchy-Schwarz inequality then gives

$$\begin{aligned}
\|\hat{\theta}_n - \theta\| &\leq O_P(1) \left\{ \int_{\delta}^{1-\delta} \|F'_{\theta}(\hat{x}(t), \hat{\theta}_n)\|^2 w(t) dt \right\}^{1/2} \\
&\quad \times \left\{ \int_{\delta}^{1-\delta} \|\hat{x}'(t) - x'_{\theta}(t)\|^2 w(t) dt \right\}^{1/2} \\
&\quad + O_P(1) \left\{ \int_{\delta}^{1-\delta} \|F'_{\theta}(\hat{x}(t), \hat{\theta}_n)\|^2 w(t) dt \right\}^{1/2}
\end{aligned}$$

$$\times \left\{ \int_{\delta}^{1-\delta} \|F(x_{\theta}(t), \theta) - F(\hat{x}(t), \theta)\|^2 w(t) dt \right\}^{1/2}.$$

By a by now standard argument, i.e. (10), (11), and the continuous mapping theorem, the righthand side can be further bounded to obtain

$$(34) \quad \|\hat{\theta}_n - \theta\| \leq O_P(1) \left\{ \sqrt{\left(b^{\alpha-1} + \frac{1}{nb^3}\right)^2 + \frac{\log n}{nb^3}} + \sqrt{\left(b^{\alpha} + \frac{1}{nb^2}\right)^2 + \frac{\log n}{nb}} \right\}.$$

Summarising the above results, we finally get that the second term at the righthand side of (31) satisfies

$$\begin{aligned} & \left\| \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta))^T (\hat{x}'(t) - x'_{\theta}(t)) w(t) dt \right\| \\ & \leq O_P(1) \sqrt{\left(b^{\alpha-1} + \frac{1}{nb^3}\right)^2 + \frac{\log n}{nb^3}} \\ & \times \left\{ \sqrt{\left(b^{\alpha-1} + \frac{1}{nb^3}\right)^2 + \frac{\log n}{nb^3}} + \sqrt{\left(b^{\alpha} + \frac{1}{nb^2}\right)^2 + \frac{\log n}{nb}} \right\} \\ & = o_P(n^{-1/2}), \end{aligned}$$

where the last equality follows from our conditions on  $b$ . Here we also see that the condition  $\alpha \geq 3$  is needed for the conclusion to hold.

To conclude the proof, it remains to consider the second term within brackets in (30). We have

$$\begin{aligned} & \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n))^T (F(x_{\theta}(t), \theta) - F(\hat{x}(t), \theta)) w(t) dt \\ (35) \quad & = \int_{\delta}^{1-\delta} (F'_{\theta}(x_{\theta}(t), \theta))^T (F(x_{\theta}(t), \theta) - F(\hat{x}(t), \theta)) w(t) dt \\ & \quad + \int_{\delta}^{1-\delta} (F'_{\theta}(\hat{x}(t), \hat{\theta}_n) - F'_{\theta}(x_{\theta}(t), \theta))^T (F(x_{\theta}(t), \theta) - F(\hat{x}(t), \theta)) w(t) dt. \end{aligned}$$

This can be analysed in a by now routine fashion, but we provide proofs. We first study the first term at the righthand side. By a standard argument we have

$$\begin{aligned} & \int_{\delta}^{1-\delta} (F'_{\theta}(x_{\theta}(t), \theta))^T (F(x_{\theta}(t), \theta) - F(\hat{x}(t), \theta)) w(t) dt \\ & = \int_{\delta}^{1-\delta} (F'_{\theta}(x_{\theta}(t), \theta))^T \int_0^1 F'_x(x_{\theta}(t) + \lambda(\hat{x}(t) - x_{\theta}(t)), \theta) d\lambda (\hat{x}(t) - x_{\theta}(t)) w(t) dt \\ & \quad = \int_{\delta}^{1-\delta} (F'_{\theta}(x_{\theta}(t), \theta))^T F'_x(x_{\theta}(t), \theta) (\hat{x}(t) - x_{\theta}(t)) w(t) dt \\ & + \int_{\delta}^{1-\delta} (F'_{\theta}(x_{\theta}(t), \theta))^T \int_0^1 [F'_x(x_{\theta}(t) + \lambda(\hat{x}(t) - x_{\theta}(t)), \theta) - F'_x(x_{\theta}(t), \theta)] d\lambda (\hat{x}(t) - x_{\theta}(t)) w(t) dt \\ & = T_3 + T_4. \end{aligned}$$

Recalling (17), we see that  $T_3$  appears in the leading term  $\Gamma(\hat{x}) - \Gamma(x_{\theta})$  in (16) and completes it together with the first term at the righthand side of (31). Next we

consider  $T_4$ . Introduce the notation  $F'_x(x, \theta) = B(x, \theta) = (b_{i,j}(x, \theta))_{i,j}$ . We have

$$\begin{aligned}
& \left\| \int_0^1 [F'_x(x_\theta(t) + \lambda(\hat{x}(t) - x_\theta(t)), \theta) - F'_x(x_\theta(t), \theta)] d\lambda (\hat{x}(t) - x_\theta(t)) \right\| \\
& \leq \left( \sup_{t \in [\delta, 1-\delta]} \|\hat{x}(t) - x_\theta(t)\| \right) \\
& \quad \times \int_0^1 \|F'_x(x_\theta(t) + \lambda(\hat{x}(t) - x_\theta(t)), \theta) - F'_x(x_\theta(t), \theta)\| d\lambda \\
& \leq \left( \sup_{t \in [\delta, 1-\delta]} \|\hat{x}(t) - x_\theta(t)\| \right) \\
& \quad \times \int_0^1 \sum_{i,j} |b_{i,j}(x_\theta(t) + \lambda(\hat{x}(t) - x_\theta(t)), \theta) - b_{i,j}(x_\theta(t), \theta)| d\lambda \\
& \leq \left( \sup_{t \in [\delta, 1-\delta]} \|\hat{x}(t) - x_\theta(t)\| \right) \\
& \quad \times \sum_{i,j} \int_0^1 \left\| \int_0^1 \frac{\partial}{\partial x} b_{i,j}(x_\theta(t) + \kappa\lambda(\hat{x}(t) - x_\theta(t)), \theta) d\kappa \lambda (\hat{x}(t) - x_\theta(t)) \right\| d\lambda \\
& \leq \left( \sup_{t \in [\delta, 1-\delta]} \|\hat{x}(t) - x_\theta(t)\|^2 \right) \\
& \quad \times \sum_{i,j} \int_0^1 \int_0^1 \left\| \frac{\partial}{\partial x} b_{i,j}(x_\theta(t) + \kappa\lambda(\hat{x}(t) - x_\theta(t)), \theta) \right\| d\kappa d\lambda,
\end{aligned}$$

where in the last inequality we used the fact that  $0 \leq \lambda \leq 1$ . Since by convergence in probability of  $\hat{x}$  to  $x_\theta$ , Condition 4 and the continuous mapping theorem the integrals on the righthand side of the above display are bounded in probability, it follows from (10) that  $\|T_4\|$  is

$$O_P(1) \left\{ \left( b^\alpha + \frac{1}{nb^2} \right)^2 + \frac{\log n}{nb} \right\}.$$

This in turn is  $o_P(n^{-1/2})$  because of the conditions on  $b$ . Finally, we treat the second term at the righthand side of (35). By the Cauchy-Schwarz inequality, its norm can be bounded by

$$\begin{aligned}
& \left\{ \int_\delta^{1-\delta} \|F'_\theta(\hat{x}(t), \hat{\theta}_n) - F'_\theta(x_\theta(t), \theta)\|^2 w(t) dt \right\}^{1/2} \\
& \quad \times \left\{ \int_\delta^{1-\delta} \|F(x_\theta(t), \theta) - F(\hat{x}(t), \theta)\|^2 w(t) dt \right\}^{1/2}.
\end{aligned}$$

Each of the terms at the righthand side have already been treated above, see (32) and (34), and it follows that the lefthand side of the last display is  $o_P(n^{-1/2})$ . This concludes the proof of Proposition 3.  $\square$



*Proof of Proposition 4.* By a standard decomposition, we have

$$\begin{aligned}\mathbb{E}[(\Delta(\hat{\mu}_n) - \Delta(\mu))^2] &= (\mathbb{E}[\Delta(\hat{\mu}_n)] - \Delta(\mu))^2 + \text{Var}[\Delta(\hat{\mu}_n)] \\ &= T_1^2 + T_2.\end{aligned}$$

The statement of the theorem will follow from Chebyshev's inequality, provided we show that the righthand side of the above display is  $O(n^{-1})$ . For  $T_1$  we have

$$\begin{aligned}|T_1| &= \left| \int_{\mathbb{R}} v(t)k(t)(\mathbb{E}[\hat{\mu}_n(t)] - \mu(t))dt \right| \\ &\leq \sup_{t \in [\delta, 1-\delta]} |\mathbb{E}[\hat{\mu}_n(t)] - \mu(t)| \int_{\mathbb{R}} |v(t)k(t)| dt \\ &= O\left(b^\alpha + \frac{1}{nb^2}\right),\end{aligned}$$

where the last equality follows from (20). Taking  $1/(2\alpha) \leq \gamma \leq 1/4$  gives that  $T_1$  is  $O(n^{-1/2})$ . We next consider  $T_2$ . By independence of the  $\epsilon_i$ 's and the fact that  $\max_i |t_i - t_{i-1}| \lesssim n^{-1}$ , we have

$$\begin{aligned}T_2 &= \text{Var} \left[ \sum_{i=1}^n (t_i - t_{i-1}) Y_i \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K\left(\frac{t-t_i}{b}\right) dt \right] \\ &\lesssim \frac{\sigma^2}{n} \sum_{i=1}^n (t_i - t_{i-1}) \left( \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K\left(\frac{t-t_i}{b}\right) dt \right)^2 \\ &= \frac{\sigma^2}{n} \int_0^1 \left( \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K\left(\frac{t-s}{b}\right) dt \right)^2 ds \\ &\quad + \frac{\sigma^2}{n} \left\{ \sum_{i=1}^n (t_i - t_{i-1}) \left( \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K\left(\frac{t-t_i}{b}\right) dt \right)^2 \right. \\ &\quad \left. - \int_0^1 \left( \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K\left(\frac{t-s}{b}\right) dt \right)^2 ds \right\} \\ &= T_3 + T_4.\end{aligned}$$

Notice that by a change of the integration variable  $(t-s)/b = u$  we have

$$\begin{aligned}\int_{\delta}^{1-\delta} \left| v(t)k(t) \frac{1}{b} K\left(\frac{t-s}{b}\right) \right| dt &= \int_{(\delta-s)/b}^{(1-\delta-s)/b} |v(s+bu)k(s+bu)K(u)| du \\ &\leq \int_{-1}^1 |v(s+bu)k(s+bu)K(u)| du \\ &\lesssim \sup_{z \in [\delta, 1-\delta]} |v(z)k(z)|,\end{aligned}$$

where we used the fact that  $0 \leq s \leq 1$ , which implies  $|s+bu| \leq 1+b$ , as well as the fact that  $K$  is supported on  $[-1, 1]$ , while  $k$  has support on  $[\delta, 1-\delta]$ . It then follows that  $T_3 = O(n^{-1})$ . To complete the proof, it remains to bound  $T_4$ . By a standard argument we get

$$|T_4| = \frac{\sigma^2}{n} \left| \sum_{i=1}^n \int_{t_{j-1}}^{t_i} \left( \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K\left(\frac{t-t_i}{b}\right) dt \right)^2 ds \right|$$

$$\begin{aligned}
& \left| - \sum_{i=1}^n \int_{t_{j-1}}^{t_i} \left( \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K \left( \frac{t-s}{b} \right) dt \right)^2 ds \right| \\
& \leq \frac{\sigma^2}{n} \sum_{i=1}^n \int_{t_{j-1}}^{t_i} \left| \left( \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K \left( \frac{t-t_i}{b} \right) dt \right)^2 \right. \\
& \quad \left. - \left( \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K \left( \frac{t-s}{b} \right) dt \right)^2 \right| ds \\
& \leq \frac{2\sigma^2 \|vk\|_{\infty} \|K\|_{\infty}}{nb} \sum_{i=1}^n \int_{t_{j-1}}^{t_i} \left| \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K \left( \frac{t-t_i}{b} \right) dt \right. \\
& \quad \left. - \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K \left( \frac{t-s}{b} \right) dt \right| ds.
\end{aligned}$$

Now notice that for  $s \in [t_{j-1}, t_i]$  by continuous differentiability of  $K$

$$\begin{aligned}
& \left| \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K \left( \frac{t-t_i}{b} \right) dt - \int_{\delta}^{1-\delta} v(t)k(t) \frac{1}{b} K \left( \frac{t-s}{b} \right) dt \right| \\
& \leq \int_{\delta}^{1-\delta} |v(t)k(t)| \left| \frac{1}{b} K \left( \frac{t-t_i}{b} \right) - \frac{1}{b} K \left( \frac{t-s}{b} \right) \right| dt \\
& \lesssim \|K'\|_{\infty} \int_{\delta}^{1-\delta} |v(t)k(t)| dt \frac{1}{nb^2}
\end{aligned}$$

holds. Hence  $T_4$  is  $O(n^{-2}b^{-3})$ . Next  $n/(n^2b^3) \rightarrow 0$  if  $\gamma < 1/3$ , and thus certainly if  $\gamma \leq 1/4$ . Consequently,  $T_4 = o(n^{-1})$ . This completes the proof of Proposition 4.  $\square$

*Proof of Theorem 1.* The result is an easy consequence of Propositions 3 and 4.  $\square$

#### APPENDIX A

The proof of Proposition 1 is based on the following two lemmas, which provide integral approximations to the bias and variance of the estimator  $\hat{\mu}_n$  and its derivative  $\hat{\mu}'_n$  at a point  $t$ .

**Lemma 1.** *Let  $\mu$  and  $K$  be continuously differentiable and let  $K$  be supported on the interval  $[-1, 1]$ . For any  $t \in [0, 1]$*

$$(36) \quad \mathbb{E}[\hat{\mu}_n(t)] = \int_0^1 \mu(s) \frac{1}{b} K \left( \frac{t-s}{b} \right) ds + O \left( \frac{1}{nb^2} \right)$$

*holds in the regression model (6). The order bound on the remainder term in (36) is uniform in  $t \in [0, 1]$ .*

*Proof.* The proof is based on the Riemann sum approximation of the integral. Since  $\mathbb{E}[\epsilon_i] = 0$ , we have

$$\begin{aligned}
\mathbb{E}[\hat{\mu}_n(t)] &= \int_0^1 \mu(s) \frac{1}{b} K \left( \frac{t-s}{b} \right) ds \\
&\quad - \int_0^1 \mu(s) \frac{1}{b} K \left( \frac{t-s}{b} \right) ds + \sum_{i=1}^n (t_i - t_{i-1}) \mu(t_i) \frac{1}{b} K \left( \frac{t-t_i}{b} \right).
\end{aligned}$$

The first term at the righthand side of this expression is the first term of (36). We will now establish an upper bound on the difference of the other two terms. Using continuous differentiability of  $\mu$  and  $K$  and the fact that  $\max_i |t_i - t_{i-1}| = O(n^{-1})$ , we have

$$\begin{aligned} & \left| \int_0^1 \mu(s) \frac{1}{b} K\left(\frac{t-s}{b}\right) ds - \sum_{i=1}^n (t_i - t_{i-1}) \mu(t_i) \frac{1}{b} K\left(\frac{t-t_i}{b}\right) \right| \\ &= \left| \sum_{i=1}^n \int_{t_{j-1}}^{t_i} \left\{ \mu(s) \frac{1}{b} K\left(\frac{t-s}{b}\right) - \mu(t_i) \frac{1}{b} K\left(\frac{t-t_i}{b}\right) \right\} ds \right| \\ &\leq \sum_{i=1}^n \int_{t_{j-1}}^{t_i} \left| \mu(s) \frac{1}{b} K\left(\frac{t-s}{b}\right) - \mu(s) \frac{1}{b} K\left(\frac{t-t_i}{b}\right) \right| ds \\ &\quad + \sum_{i=1}^n \int_{t_{j-1}}^{t_i} \left| \mu(s) \frac{1}{b} K\left(\frac{t-t_i}{b}\right) - \mu(t_i) \frac{1}{b} K\left(\frac{t-t_i}{b}\right) \right| ds \\ &\lesssim \frac{1}{nb^2} \|\mu\|_\infty \|K'\|_\infty + \frac{1}{nb} \|\mu'\|_\infty \|K\|_\infty, \end{aligned}$$

which is of order  $n^{-1}b^{-2}$ . This establishes (36).  $\square$

The second lemma can be proved along the same lines as the previous one and therefore we omit its proof. The existence of the second derivative of  $K$  is needed in the proof of this lemma.

**Lemma 2.** *Let  $\mu$  be continuously differentiable and let  $K$  be twice continuously differentiable and be supported on the interval  $[-1, 1]$ . For all  $t \in [0, 1]$*

$$(37) \quad \mathbb{E}[\hat{\mu}'_n(t)] = \int_0^1 \mu(s) \frac{1}{b^2} K'\left(\frac{t-s}{b}\right) ds + O\left(\frac{1}{nb^3}\right)$$

*holds in the regression model (6). Furthermore, if  $b \leq \delta$  and  $t \in [\delta, 1 - \delta]$ , then integration by parts yields*

$$(38) \quad \mathbb{E}[\hat{\mu}'_n(t)] = \int_{-1}^1 \mu'(t - bu) K(u) du + O\left(\frac{1}{nb^3}\right).$$

*The order bounds on the remainder terms in (37) and (38) are uniform in  $t$ .*

The following lemma is used in the proof of Proposition 2.

**Lemma 3.** *Let the stochastic process  $X = (X_{n,\eta})_{\eta \in \Theta}$  be defined as*

$$X = (X_{n,\eta})_{\eta \in \Theta} = \left( \int_{\delta}^{1-\delta} \|F(\hat{x}(t), \eta) - F(x_\theta(t), \eta)\|^2 w(t) dt \right)_{\eta \in \Theta}.$$

*Then under the conditions of Proposition 2 we have  $X \xrightarrow{P} 0$ , where 0 at the righthand side denotes the zero process on  $\Theta$  and convergence is understood as convergence for random elements with values in the space  $C(\Theta)$  of continuous functions on  $\Theta$ , which is equipped with the supremum norm.*

*Proof.* To prove the lemma, we will verify the conditions of Theorem 18.14 of van der Vaart (1998). By (10) and the continuous mapping theorem, see Theorem

18.11 in van der Vaart (1998), for every fixed  $\eta$  it holds that

$$(39) \quad \int_{\delta}^{1-\delta} \|F(\hat{x}(t), \eta) - F(x_{\theta}(t), \eta)\|^2 w(t) dt \xrightarrow{P} 0.$$

Consequently, for any positive integer  $k$  and any  $\eta_1, \dots, \eta_k \in \Theta$  we have

$$(X_{n,\eta_1}, \dots, X_{n,\eta_k}) \rightsquigarrow \underbrace{(0, \dots, 0)}_k$$

and hence condition (i) of Theorem 18.14 in van der Vaart (1998) is satisfied. Introduce

$$G = \bigcap_{j=1}^d \left\{ \sup_{t \in [\delta, 1-\delta]} |\hat{x}_j(t) - x_{\theta_j}(t)| \leq \beta \right\}$$

and notice

$$G^c = \bigcup_{j=1}^d \left\{ \sup_{t \in [\delta, 1-\delta]} |\hat{x}_j(t) - x_{\theta_j}(t)| > \beta \right\}.$$

For any positive  $\varepsilon$  and  $\beta$  and any partition  $\Theta_1, \dots, \Theta_m$  of  $\Theta$  we have

$$(40) \quad \begin{aligned} & P \left( \sup_{\ell} \sup_{\eta, \zeta \in \Theta_{\ell}} |X_{n,\eta} - X_{n,\zeta}| \geq \varepsilon \right) \\ & \leq P \left( \sup_{\ell} \sup_{\eta, \zeta \in \Theta_{\ell}} |X_{n,\eta} - X_{n,\zeta}| \geq \varepsilon; G \right) + P(G^c). \end{aligned}$$

By (10) we know that

$$(41) \quad \lim_{n \rightarrow \infty} P(G^c) \leq \lim_{n \rightarrow \infty} \sum_{j=1}^d P \left( \sup_{t \in [\delta, 1-\delta]} |\hat{x}_j(t) - x_{\theta_j}(t)| > \beta \right) = 0.$$

We will now show that for arbitrarily small positive  $\rho$  and  $\varepsilon$  there exists a partition  $\Theta_1, \dots, \Theta_m$  of  $\Theta$ , such that

$$\limsup_{n \rightarrow \infty} P \left( \sup_{\ell} \sup_{\eta, \zeta \in \Theta_{\ell}} |X_{n,\eta} - X_{n,\zeta}| \geq \varepsilon; G \right) \leq \rho.$$

Together with (40) and (41) this will imply condition (ii) of Theorem 18.14 in van der Vaart (1998) and hence also the fact that  $X$  converges weakly to zero. The statement of the lemma will then be a simple consequence of the fact that convergence in distribution and in probability are equivalent for constants, see Theorem 18.10 of van der Vaart (1998).

Notice that

$$\begin{aligned} & |X_{n,\eta} - X_{n,\zeta}| \\ & \leq \int_{\delta}^{1-\delta} \|F(\hat{x}(t), \eta) - F(x_{\theta}(t), \eta) - F(\hat{x}(t), \zeta) + F(x_{\theta}(t), \zeta)\| \\ & \quad \times (\|F(\hat{x}(t), \eta) - F(x_{\theta}(t), \eta)\| + \|F(\hat{x}(t), \zeta) - F(x_{\theta}(t), \zeta)\|) w(t) dt \\ & \leq \left\{ \int_{\delta}^{1-\delta} \|F(\hat{x}(t), \eta) - F(x_{\theta}(t), \eta) - F(\hat{x}(t), \zeta) + F(x_{\theta}(t), \zeta)\|^2 w(t) dt \right\}^{1/2} \end{aligned}$$

$$\begin{aligned} & \times \left\{ \int_{\delta}^{1-\delta} (\|F(\hat{x}(t), \eta) - F(x_{\theta}(t), \eta)\| + \|F(\hat{x}(t), \zeta) - F(x_{\theta}(t), \zeta)\|)^2 w(t) dt \right\}^{1/2} \\ & = \sqrt{T_3} \sqrt{T_4}. \end{aligned}$$

For  $T_3$  we have

$$\begin{aligned} T_3 \leq 2 \int_{\delta}^{1-\delta} \|F(\hat{x}(t), \eta) - F(\hat{x}(t), \zeta)\|^2 w(t) dt \\ + 2 \int_{\delta}^{1-\delta} \|F(x_{\theta}(t), \eta) - F(x_{\theta}(t), \zeta)\|^2 w(t) dt. \end{aligned}$$

Restricting  $\omega$ 's from the sample space  $\Omega$  to the set  $G$ , we get by Taylor

$$\begin{aligned} T_3 & \leq 2 \int_{\delta}^{1-\delta} \int_0^1 \|F'_{\theta}(\hat{x}(t), \zeta + \lambda(\eta - \zeta))\|^2 d\lambda \|\eta - \zeta\|^2 w(t) dt \\ & \quad + 2 \int_{\delta}^{1-\delta} \int_0^1 \|F'_{\theta}(x_{\theta}(t), \zeta + \lambda(\eta - \zeta))\|^2 d\lambda \|\eta - \zeta\|^2 w(t) dt \\ & \leq 4 \|\eta - \zeta\|^2 \int_{\delta}^{1-\delta} w(t) dt \sup_{\substack{\|x_j\| \leq \|x_{\theta_j}\|_{\infty} + \beta, j=1, \dots, d \\ \nu \in \Theta}} \|F'_{\theta}(x, \nu)\| = C(\beta, w, \theta, \Theta) \|\eta - \zeta\|^2 \end{aligned}$$

on the set  $G$ . Notice that  $C(\beta, w, \theta, \Theta)$  is a finite constant, because  $\|F'_{\theta}(x, \nu)\|$  is continuous and its supremum is taken over a compact set. By similar techniques one can show that  $T_4 \leq C'(\beta, w, \theta, \Theta)$  for some constant  $C'(\beta, w, \theta, \Theta)$  which depends only on  $\beta, w, \theta$ , and  $\Theta$ . Consequently,

$$\begin{aligned} (42) \quad & P \left( \sup_{\ell} \sup_{\eta, \zeta \in \Theta_{\ell}} |X_{n, \eta} - X_{n, \zeta}| \geq \varepsilon; G \right) \\ & \leq P \left( \sup_{\ell} \sup_{\eta, \zeta \in \Theta_{\ell}} \sqrt{C(\beta, w, \theta, \Theta) C'(\beta, w, \theta, \Theta)} \|\eta - \zeta\| \geq \varepsilon \right). \end{aligned}$$

Now take a partition  $\Theta_1, \dots, \Theta_m$  of  $\Theta$  such that for all  $\ell = 1, \dots, m$

$$0 < \text{diam } \Theta_{\ell} < \frac{\varepsilon}{\sqrt{C(\beta, w, \theta, \Theta) C'(\beta, w, \theta, \Theta)}}$$

holds, where  $\text{diam } \Theta_{\ell}$  denotes the diameter of the set  $\Theta_{\ell}$ . Observe that since  $\Theta \subset \mathbb{R}^p$  is compact, there indeed exists a finite  $m$  for which this is satisfied. The righthand side of (42) for such a partition is zero and consequently the conditions (i) and (ii) of Theorem 18.14 of van der Vaart (1998) hold. This completes the proof of the lemma.  $\square$

In a similar fashion one can prove the following lemma, which is also used in the proof of Proposition 2. We omit the proof.

**Lemma 4.** *Let the stochastic process  $X = (X_{n, \eta})_{\eta \in \Theta}$  be defined as*

$$X = (X_{n, \eta})_{\eta \in \Theta} = \left( \int_{\delta}^{1-\delta} \|\hat{x}'(t) - F(\hat{x}(t), \eta)\|^2 w(t) dt \right)_{\eta \in \Theta}.$$

*Then under the conditions of Proposition 2 we have  $X \xrightarrow{P} 0$ , where 0 at the righthand side denotes the zero process on  $\Theta$  and convergence is understood as convergence*

for random elements with values in the space  $C(\Theta)$  of continuous functions on  $\Theta$ , which is equipped with the supremum norm.

## APPENDIX B

Here we state and prove a modification of Proposition 1 for the case when the  $\epsilon_i$ 's are bounded.

**Proposition 5.** *In the regression model (6) replace the assumption of Gaussianity of the  $\epsilon_i$ 's by  $|\epsilon_i| \leq C$  for some constant  $C > 0$  and suppose Condition 5 holds.*

(i) *If  $\mu$  is  $\alpha \geq 1$  times continuously differentiable and  $b \rightarrow 0$  as  $n \rightarrow \infty$ , then*

$$(43) \quad \sup_{t \in [\delta, 1-\delta]} |\hat{\mu}_n(t) - \mu(t)| = O_P \left( \sqrt{\left(b^\alpha + \frac{1}{nb^2}\right)^2 + \frac{\log n}{nb}} \right).$$

(ii) *If  $\mu$  is  $\alpha \geq 2$  times continuously differentiable and  $b \rightarrow 0$  as  $n \rightarrow \infty$ , then*

$$(44) \quad \sup_{t \in [\delta, 1-\delta]} |\hat{\mu}'_n(t) - \mu'(t)| = O_P \left( \sqrt{\left(b^{\alpha-1} + \frac{1}{nb^3}\right)^2 + \frac{\log n}{nb^3}} \right)$$

is valid. Moreover,  $\hat{\mu}_n$  and  $\hat{\mu}'_n$  are consistent on  $[\delta, 1-\delta]$ , if  $nb^3/\log n \rightarrow \infty$  holds additionally.

*Proof.* The proof of (43) follows the same steps as the proof of (8). The only difference is that we need to show that

$$(45) \quad \mathbb{E} \left[ \max_{1 \leq j \leq N} |Z_j|^2 \right] = O \left( \frac{\log n}{nb} \right)$$

holds also for bounded  $\epsilon_i$ 's and not only for the Gaussian  $\epsilon_i$ 's. To this end we will use some results from Chapter 2.2 of Wellner and van der Vaart (1996). Let  $\eta$  be a nondecreasing and convex function on  $[0, \infty)$ , such that  $\eta(0) = 0$ . The Orlicz norm  $\|X\|_\eta$  of a random variable  $X$  is defined as

$$\|X\|_\eta = \inf \left\{ C > 0 : \mathbb{E} \left[ \eta \left( \frac{|X|}{C} \right) \right] \leq 1 \right\}.$$

A particular  $\eta$  that we will use is  $\eta_p(x) = e^{x^p} - 1$  for  $p \geq 1$ . Since the  $\epsilon_i$ 's have mean zero and are bounded, for any  $x > 0$  Hoeffding's inequality, see Hoeffding (1963), implies

$$P(|Z_j| > x) \leq 2 \exp \left( -2x^2 / \left( \sum_{i=1}^n C^2 (S_i(s_j))^2 \right) \right).$$

By Condition 1 it holds that

$$\begin{aligned} C^2 \sum_{i=1}^n (S_i(s_j))^2 &\lesssim C^2 \|K\|_\infty^2 \frac{1}{n^2 b^2} \sum_{i=1}^n \mathbf{1}_{\|s_j - t_i\| \leq b} \\ &\leq \frac{1}{nb} C^2 \|K\|_\infty^2 c_1 \max \left( 2, \max_n \frac{1}{nb} \right) = \frac{1}{C_0 nb}. \end{aligned}$$

Thus the inequality

$$P(|Z_j| > x) \leq 2 \exp(-2C_0 nbx^2)$$

is valid. By Lemma 2.2.1 of Wellner and van der Vaart (1996) it then follows that

$$(46) \quad \max_j \|Z_j\|_{\eta_2} \leq \frac{C_1}{\sqrt{nb}},$$

where  $C_1$  depends on  $C_0$  only. Let  $\|X\|_2$  denote the  $L_2$  norm of a random variable  $X$ , i.e.  $\|X\|_2 = \sqrt{\mathbb{E}[X^2]}$ . Notice that the inequality

$$(47) \quad \|X\|_2 \leq \|X\|_{\eta_2},$$

holds, because of  $e^{x^2} - 1 \geq x^2$ . The inequalities (46) and (47) combined with Lemma 2.2.2 of Wellner and van der Vaart (1996) yield that

$$\sqrt{\mathbb{E} \left[ \max_{1 \leq j \leq N} |Z_j|^2 \right]} \leq \frac{C_3}{\sqrt{nb}} \eta_2^{-1}(N),$$

where the constant  $C_3$  is independent of  $N$ . Now notice that for  $N \geq 4$

$$\eta_2^{-1}(N) = \sqrt{\log(N+1)} \leq \sqrt{\log(N^2)} = 2\sqrt{\log n}.$$

Hence (45) holds and this completes the proof of (43). Formula (44) can be proved in a similar fashion.  $\square$

#### REFERENCES

- V.I. Arnold. *Ordinary Differential Equations*. MIT Press, Massachusetts, 1973.
- J.K. Benedetti. On the nonparametric estimation of regression functions. *J. Roy. Statist. Soc. Ser. B*, 39:248–253, 1977.
- P.J. Bickel and Y. Ritov. Nonparametric estimators which can be “plugged-in”. *Ann. Statist.*, 31:1033–1053, 2003.
- P.J. Bickel, C.A.J. Klaassen, Y. Ritov and J.A. Wellner. *Efficient and Adaptive Estimation for Semiparametric Models*. Springer, New York, 1998.
- N.J-B. Brunel. Parameter estimation of ODE's via nonparametric estimators. *Electron. J. Stat.*, 2:1242–1267, 2008.
- I-C. Chou and E.O. Voit. Recent developments in parameter estimation and structure identification of biochemical and genomic systems. *Math. Biosci.*, 219:57–83, 2009.
- L. Edelstein-Keshet. *Mathematical Models in Biology*. Society for Industrial and Applied Mathematics, Philadelphia, 2005.
- A.J. van Es. *Aspects of Nonparametric Density Estimation*. CWI Tract, 77. Stichting Mathematisch Centrum, Centrum voor Wiskunde en Informatica, Amsterdam, 1991.
- M. Feinberg. *Lectures on Chemical Reaction Networks*. Lectures delivered at the Mathematics Research Center, University of Wisconsin-Madison, 1979.
- Th. Gasser, H-G. Müller and V. Mammitzsch. Kernels for nonparametric curve estimation. *J. Roy. Statist. Soc. Ser. B*, 47:238–252, 1985.
- Th. Gasser and H-G. Müller. Estimating regression functions and their derivatives by the kernel method. *Scand. J. Statist.*, 11:197–211, 1984.
- S. van de Geer. Estimating a regression function. *Ann. Statist.*, 18:907-924, 1990.
- S. van de Geer and M. Wegkamp. Consistency for the least squares estimator in nonparametric regression. *Ann. Statist.*, 24:2513-2523, 1996.
- A. Gelman, F.Y. Bois and J. Jiang. Physiological pharmacokinetic analysis using population modeling and informative prior distributions. *J. Amer. Statist. Assoc.*, 91:1400–1412, 1996.

- M. Girolami. Bayesian inference for differential equations. *Theor. Comput. Sci.*, 408:4–16, 2008.
- L. Goldstein and K. Messer. Optimal plug-in estimators for nonparametric functional estimation. *Ann. Statist.*, 20:1306–1328, 1992.
- E. Hairer and G. Wanner. *Solving Ordinary Differential Equations. II. Stiff and Differential-Algebraic Problems*. Springer-Verlag, Berlin, 1996.
- P. Hall and J.S. Marron. On variance estimation in nonparametric regression. *Biometrika*, 77:415–419, 1990.
- P.W. Hemker. Numerical methods for differential equations in system simulation and in parameter estimation. In H.C. Hemker and B. Hess, editors, *Analysis and Simulation of Biochemical Systems*, North Holland Publ. Comp., Amsterdam, 59–80, 1972.
- W.S. Hlavacek and M.A. Savageau. Rules for coupled expression of regulator and effector genes in inducible circuits. *J. Mol. Biol.*, 255:121–139, 1996.
- W. Hoeffding. Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.*, 58:13–30, 1963.
- P.J. Huber. *Robust Statistics*. John Wiley & Sons, Inc., New York, 1981.
- R.I. Jennrich. Asymptotic properties of non-linear least squares estimators. *Ann. Math. Statist.*, 40:633–643, 1969.
- M.C. Jones, J.S. Marron and S.J. Sheather. A brief survey of bandwidth selection for density estimation. *J. Amer. Statist. Assoc.*, 91:401–407, 1996.
- S. Kikuchi, D. Tominaga, M. Arita, K. Takahashi and M. Tomita. Dynamic modeling of genetic networks using genetic algorithm and S-system. *Bioinformatics*, 19:643–650, 2003.
- D.W. Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *SIAM J. Appl. Math.*, 11:431–441, 1963.
- H. Liang and H. Wu. Parameter estimation for differential equation models using a framework of measurement error in regression models. *J. Amer. Statist. Assoc.*, 103:1570–1583, 2008.
- C.R. Loader. Bandwidth selection: classical or plug-in? *Ann. Statist.*, 27:415–438, 1999.
- K. Messer and L. Goldstein. A new class of kernels for nonparametric curve estimation. *Ann. Statist.*, 21:179–195, 1993.
- D. Pollard and P. Radchenko. Nonlinear least-squares estimation. *J. Multivariate Anal.*, 97:548–562, 2006.
- M.B. Priestley and M.T. Chao. Non-parametric function fitting. *J. Roy. Statist. Soc. Ser. B*, 34:385–392, 1972.
- X. Qi and H. Zhao. Asymptotic efficiency and finite sample properties of the generalized profiling estimation of parameters in ordinary differential equations. *Ann. Statist.*, 38:435–481, 2010.
- J.O. Ramsay, G. Hooker, D. Campbell and J. Cao. Parameter estimation for differential equations: a generalized smoothing approach. With discussions and a reply by the authors. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 69:741–796, 2007.
- E. Schuster and S. Yakowitz. Contributions to the theory of nonparametric regression, with application to system identification. *Ann. Statist.*, 7:139–149, 1979.
- E.D. Sontag. Structure and stability of certain chemical networks and applications to the kinetic proofreading model of T-cell receptor signal transduction. *IEEE Trans. Automat. Control*, 46:1028–1047, 2001.



- W.J.H. Stortelder. Parameter estimation in dynamic systems. *Math. Comput. Simulat.*, 42:135–142, 1996.
- A.B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, New York, 2009.
- A.W. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, Cambridge, 1998.
- A.W. van der Vaart and J.A. Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer-Verlag, New York, 1996.
- J.M. Varah. A spline least squares method for numerical parameter estimation in differential equations. *SIAM J. Sci. Stat. Comput.*, 3:28–46, 1982.
- E.O. Voit. *Computational Analysis of Biochemical Systems: A Practical Guide for Biochemists and Molecular Biologists*. Cambridge University Press, Cambridge, 2000.
- E.O. Voit and J. Almeida. Decoupling dynamical systems for pathway identification from metabolic profiles. *Bioinformatics*, 10:1670–1681, 2004.
- E.O. Voit and M.A. Savageau. Power-law approach to modeling biological systems; III. Methods of analysis. *J. Ferment. Technol.*, 60:233–241, 1982.
- M.P. Wand and M.C. Jones. *Kernel Smoothing*. Chapman and Hall, London, 1995.
- L. Wasserman. *All of Nonparametric Statistics*. Springer, New York, 2006.
- C-F. Wu. Asymptotic theory of nonlinear least squares estimation. *Ann. Statist.*, 9:501–513, 1981.

DEPARTMENT OF MATHEMATICS, VU UNIVERSITY AMSTERDAM, DE BOELELAAN 1081, 1081 HV AMSTERDAM, THE NETHERLANDS

*E-mail address:* `shota@few.vu.nl`

KORTEWEG-DE VRIES INSTITUTE FOR MATHEMATICS, UNIVERSITEIT VAN AMSTERDAM, P.O. BOX 94248, 1090 GE AMSTERDAM, THE NETHERLANDS

*E-mail address:* `C.A.J.Klaassen@uva.nl`