

# Supplement to “Classification of non-parametric regression functions in longitudinal data models”

Michael Vogt  
University of Bonn

Oliver Linton  
University of Cambridge

In this supplement, we investigate the finite sample performance of the bandwidth selection rule from Section 4.3 by means of a simulation study. Moreover, we provide the proofs of Theorems 1–3 that are omitted in the paper.

## 1 Additional simulations

We now investigate the performance of the bandwidth selection procedure proposed in Section 4.3 of the paper. To do so, we pick one of the clusters from our simulation setup of Section 5 and simulate data from this cluster. In particular, we consider the cluster  $G_5$  with  $n_5 = |G_5| = 10$  and  $g_5(x) = 1.75 \arctan(5(x - 0.6)) + 0.75$  and draw data from the model equation

$$Y_{it} = g_5(X_{it}) + \varepsilon_{it} \quad (1 \leq i \leq n_5, 1 \leq t \leq T), \quad (\text{S.1})$$

where the model variables  $X_{it}$  and  $\varepsilon_{it}$  are generated in exactly the same way as in the simulations.

As discussed in Section 4.3 of the paper, our bandwidth selection procedure is based on minimizing the residual sum of squares criterion  $\text{RSS}_i^{(j)}(h)$  for different pairs of indices  $(i, j)$ . More precisely, we define our bandwidth selector by

$$\hat{h} = \frac{1}{L} \sum_{1 \leq \ell \leq L} \hat{h}_{i_{2\ell-1}}^{(i_{2\ell})},$$

where  $L = n_5/2$  and  $\hat{h}_i^{(j)} = \arg\min_h \text{RSS}_i^{(j)}(h)$ . As already mentioned in Section 4.3,  $\hat{h}$  can be regarded as an approximation to the optimal bandwidth  $h^*$  in a mean integrated squared error sense, which is defined as  $h^* = \arg\min_h \text{MISE}_i(h)$ . Note that under the conditions of Section 4.3,  $\text{MISE}_i(h)$  is the same for all  $1 \leq i \leq n_5$  and thus  $h^*$  is a group-wide optimal bandwidth independent of  $i$ .

To examine the finite sample behaviour of the bandwidth estimator  $\hat{h}$ , we draw  $N = 1000$  samples from the setting (S.1) for each time series length  $T \in \{100, 150, 200, 500\}$  and compute the bandwidth  $\hat{h}$  for each simulated sample. To do so, we define an equidistant grid  $\mathcal{G}$  of step length 0.01 which spans the interval  $[0.025, 0.5]$  and minimize the criterion functions  $\text{RSS}_i^{(j)}(h)$  over all bandwidth values  $h \in \mathcal{G}$ . The optimal bandwidth  $h^*$  can be calculated to be approximately 0.225, 0.205, 0.195, 0.165 for the time series lengths  $T = 100, 150, 200, 500$ , respectively.

Figure 1 summarizes the simulation results. Each panel shows the distribution of the differences  $h^* - \hat{h}$  for a specific time series length  $T$ . In particular, the bars in the plots give the number of simulations (out of total of 1000) in which the difference  $h^* - \hat{h}$  takes a certain value. The plots of Figure 1 suggest that  $\hat{h}$  approximates the optimal value  $h^*$  reasonably well. They also make visible that the precision of the estimator  $\hat{h}$  improves quite slowly as the sample size grows. This is not surprising as the convergence rate of standard bandwidth selectors (based on cross-validation or penalization techniques) is known to be very slow; see e.g. Härdle et al. (1988).

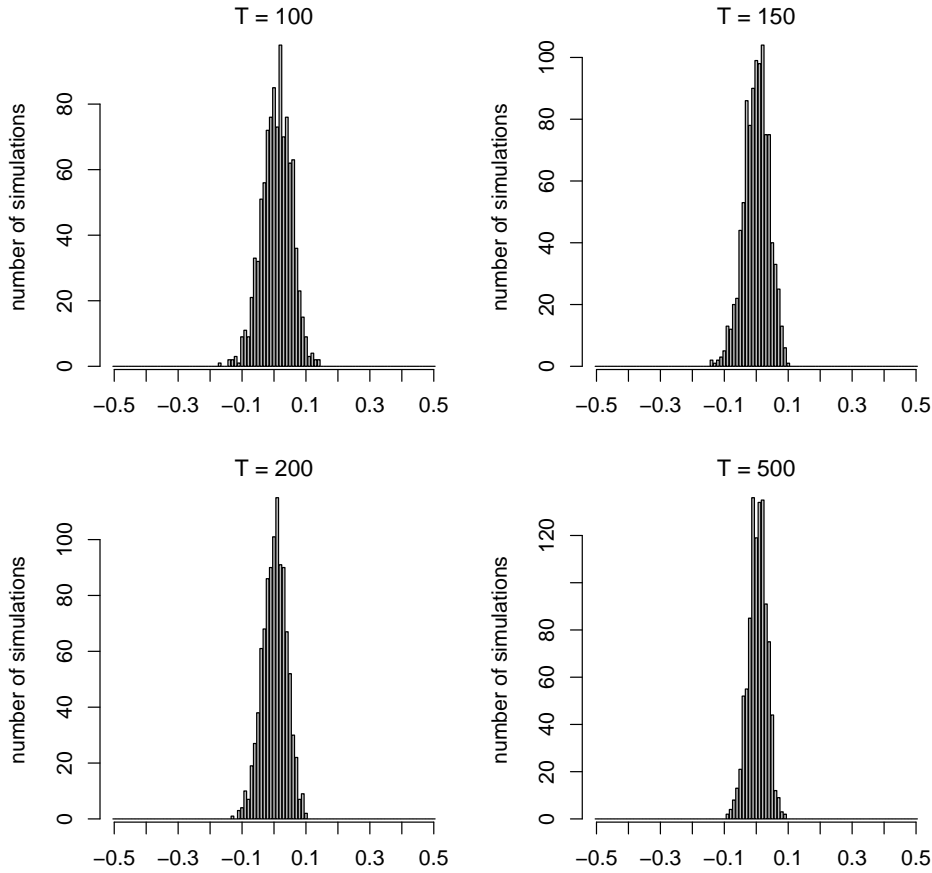


Figure 1: Simulation results for the bandwidth selection procedure from Section 4.3. Each panel depicts the distribution of the differences  $h^* - \hat{h}$  for a specific time series length  $T$ . The optimal bandwidth  $h^*$  is approximately 0.225, 0.205, 0.195, 0.165 for  $T = 100, 150, 200, 500$ , respectively.

## 2 Technical details

In this section, we prove Theorems 1–3. Throughout the section, the symbol  $C$  denotes a universal real constant which may take a different value on each occurrence.

### Auxiliary results

In the proofs of Theorems 1–3, we make use of the following uniform convergence result.

**Lemma S.1.** *Let Conditions 1–5 be satisfied, define  $I_h = [C_1h, 1 - C_1h]$  and set  $a_{n,T} = T^{-1/10}$ . It holds that*

$$\begin{aligned} \max_{1 \leq i \leq n} \sup_{x \in I_h} |\widehat{m}_i(x) - m_i(x)| &= O_p(a_{n,T} + h^2) \\ \max_{1 \leq i \leq n} \sup_{x \in [0,1] \setminus I_h} |\widehat{m}_i(x) - m_i(x)| &= O_p(a_{n,T} + h). \end{aligned}$$

If we strengthen the moment assumptions in Condition 3 to hold for some  $\theta > 20/3$ , we can improve this result to hold with  $a_{n,T} = \sqrt{\log T / (Th)}$ . From Lemma S.1, it easily follows that

$$\max_{1 \leq i, j \leq n} |\widehat{\Delta}_{ij} - \Delta_{ij}| = o_p(1). \quad (\text{S.2})$$

Moreover, we obtain that

$$\max_{i,j \in G_k} \widehat{\Delta}_{ij} = O_p(a_{n,T}^2 + h^3) \quad (\text{S.3})$$

for any  $1 \leq k \leq K$ . Notably, (S.3) merely provides an upper bound on the rate of  $\max_{i,j \in G_k} \widehat{\Delta}_{ij}$ . The reason is as follows: Directly applying Lemma S.1 does not take into account that the argument  $x$  of the smoothers  $\widehat{m}_i(x)$  and  $\widehat{m}_j(x)$  is integrated out in  $\widehat{\Delta}_{ij}$ . We now derive the sharp rate of  $\max_{i,j \in G_k} \widehat{\Delta}_{ij}$  under stronger assumptions than Conditions 1–5.

**Lemma S.2.** *Let Conditions 1–5 be satisfied, let  $h \leq CT^{-(2/9+\delta)}$  for some constant  $C$  and some small  $\delta > 0$ , and choose the weight function  $\pi$  such that its support is contained in  $I_h = [C_1h, 1 - C_1h]$  for sufficiently large  $T$ . In addition, drop the fixed effects  $\alpha_i$  and  $\gamma_t$  from the model and suppose that the following conditions hold:*

- (i) *The variables  $X_{it}$  and  $\varepsilon_{it}$  are independent both across  $i$  and  $t$ . Moreover,  $X_{it}$  and  $\varepsilon_{it}$  are independent of each other for any  $i$  and  $t$ .*
- (ii) *The second derivatives  $m_i''$  fulfill the Lipschitz condition that  $|m_i''(x) - m_i''(x')| \leq L|x - x'|$  for all  $x, x'$  and a constant  $L$  independent of  $i$ .*
- (iii) *There exist constants  $M, \gamma > 0$  such that for all indices  $i, t$  and for all  $c \geq 0$ ,  $\mathbb{P}(|\varepsilon_{it}| \geq c) \leq M \int_c^\infty \exp(-\gamma r^2) dr$ .*

Then for any  $1 \leq k \leq K$ ,

$$\max_{i,j \in G_k} \hat{\Delta}_{ij} = \max_{\substack{i,j \in G_k \\ i < j}} \frac{\mathcal{B}_{ij}}{Th} + O_p\left(\frac{\log T}{Th^{1/2}}\right) = O_p\left(\frac{1}{Th}\right),$$

where  $\mathcal{B}_{ij}$  is defined in Subsection 4.1 of the paper.

We now prove Lemmas S.1 and S.2.

**Proof of Lemma S.1.** For the proof, we modify standard arguments to derive uniform convergence rates for kernel estimators, which can be found e.g. in Masry (1996), Bosq (1998) or Hansen (2008). These arguments are designed to derive the rate of  $\sup_x |\hat{m}_i(x) - m_i(x)|$  for a fixed individual  $i$ . They thus yield the rate which is uniform over  $x$  but pointwise in  $i$ . In contrast to this, we aim to derive the rate which is uniform both over  $x$  and  $i$ . To do so, we write

$$\hat{m}_i(x) - m_i(x) = [Q_{i,V}(x) + Q_{i,B}(x) - Q_{i,\gamma}(x)] / \hat{f}_i(x) - \bar{Q}_i + \bar{\bar{Q}}_i,$$

where

$$\begin{aligned} Q_{i,V}(x) &= \frac{1}{T} \sum_{t=1}^T W_h(X_{it} - x) \varepsilon_{it} \\ Q_{i,B}(x) &= \frac{1}{T} \sum_{t=1}^T W_h(X_{it} - x) [m_i(X_{it}) - m_i(x)] \\ Q_{i,\gamma}(x) &= \frac{1}{T} \sum_{t=1}^T W_h(X_{it} - x) \left( \frac{1}{n-1} \sum_{\substack{j=1 \\ j \neq i}}^n [m_j(X_{jt}) + \varepsilon_{jt}] \right) \\ \bar{Q}_i &= \frac{1}{T} \sum_{t=1}^T [m_i(X_{it}) + \varepsilon_{it}] \\ \bar{\bar{Q}}_i &= \frac{1}{(n-1)T} \sum_{\substack{j=1 \\ j \neq i}}^n \sum_{t=1}^T [m_j(X_{jt}) + \varepsilon_{jt}] \end{aligned}$$

and  $\hat{f}_i(x) = T^{-1} \sum_{t=1}^T W_h(X_{it} - x)$ . In what follows, we show that

$$\max_{1 \leq i \leq n} \sup_{x \in [0,1]} |Q_{i,V}(x)| = O_p(a_{n,T}) \quad (\text{S.4})$$

$$\max_{1 \leq i \leq n} \sup_{x \in [0,1]} |Q_{i,B}(x) - \mathbb{E}[Q_{i,B}(x)]| = O_p(a_{n,T}) \quad (\text{S.5})$$

$$\max_{1 \leq i \leq n} \sup_{x \in [0,1]} |Q_{i,\gamma}(x)| = O_p(a_{n,T}) \quad (\text{S.6})$$

$$\max_{1 \leq i \leq n} \sup_{x \in [0,1]} |\hat{f}_i(x) - \mathbb{E}[\hat{f}_i(x)]| = O_p\left(\sqrt{\frac{\log T}{Th}}\right). \quad (\text{S.7})$$

Moreover, standard bias calculations yield that  $\max_{1 \leq i \leq n} \sup_{x \in I_h} |\mathbb{E}[Q_{i,B}(x)]| = O(h^2)$ ,  $\max_{1 \leq i \leq n} \sup_{x \in [0,1] \setminus I_h} |\mathbb{E}[Q_{i,B}(x)]| = O(h)$  and  $|\mathbb{E}[\widehat{f}_i(x)]| \geq C > 0$  for all  $x \in [0,1]$  and  $1 \leq i \leq n$  with some constant  $C$  independent of  $i$  and  $x$ . Finally, a simplified version of the arguments for (S.4) shows that  $\max_{1 \leq i \leq n} |\overline{Q}_i| = O_p(a_{n,T})$  as well as  $\max_{1 \leq i \leq n} |\overline{\overline{Q}}_i| = O_p(a_{n,T})$ . Lemma S.1 immediately follows upon combining (S.4)–(S.7) with these statements.  $\square$

**Proof of (S.4).** Set  $\psi_{n,T} = (nT)^{1/(\theta-\delta)}$ , where  $\theta$  is introduced in Condition 3 and  $\delta > 0$  is a small positive number. Moreover, define

$$\begin{aligned}\varepsilon_{it}^{\leq} &= \varepsilon_{it} 1(|\varepsilon_{it}| \leq \psi_{n,T}) \\ \varepsilon_{it}^{\geq} &= \varepsilon_{it} 1(|\varepsilon_{it}| > \psi_{n,T}).\end{aligned}$$

With this notation at hand, we can rewrite the term  $Q_{i,V}(x)$  as

$$Q_{i,V}(x) = \sum_{t=1}^T Z_{it,T}^{\leq}(x) + \sum_{t=1}^T Z_{it,T}^{\geq}(x),$$

where

$$\begin{aligned}Z_{it,T}^{\leq}(x) &= (W_h(X_{it} - x) \varepsilon_{it}^{\leq} - \mathbb{E}[W_h(X_{it} - x) \varepsilon_{it}^{\leq}]) / T \\ Z_{it,T}^{\geq}(x) &= (W_h(X_{it} - x) \varepsilon_{it}^{\geq} - \mathbb{E}[W_h(X_{it} - x) \varepsilon_{it}^{\geq}]) / T.\end{aligned}$$

We thus split  $Q_{i,V}(x)$  into the “interior part”  $\sum_{t=1}^T Z_{it,T}^{\leq}(x)$  and the “tail part”  $\sum_{t=1}^T Z_{it,T}^{\geq}(x)$ . This parallels the standard arguments for deriving the convergence rate of  $\sup_{x \in [0,1]} |Q_{i,V}(x)|$  for a fixed individual  $i$ . As we maximize over  $i$ , we however choose the truncation sequence  $\psi_{n,T}$  to go to infinity much faster than in the standard case with a fixed  $i$ .

We now proceed in several steps. To start with, we show that

$$\max_{1 \leq i \leq n} \sup_{x \in [0,1]} \left| \sum_{t=1}^T Z_{it,T}^{\geq}(x) \right| = O_p(a_{n,T}). \quad (\text{S.8})$$

This can be achieved as follows:

$$\begin{aligned}& \mathbb{P} \left( \max_{1 \leq i \leq n} \sup_{x \in [0,1]} \left| \sum_{t=1}^T Z_{it,T}^{\geq}(x) \right| > a_{n,T} \right) \\ & \leq \sum_{i=1}^n \mathbb{P} \left( \sup_{x \in [0,1]} \left| \frac{1}{T} \sum_{t=1}^T W_h(X_{it} - x) \varepsilon_{it}^{\geq} \right| > \frac{a_{n,T}}{2} \right) \\ & \quad + \sum_{i=1}^n \mathbb{P} \left( \sup_{x \in [0,1]} \left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[W_h(X_{it} - x) \varepsilon_{it}^{\geq}] \right| > \frac{a_{n,T}}{2} \right).\end{aligned}$$

With the help of Condition 3, we obtain that

$$\begin{aligned} & \sum_{i=1}^n \mathbb{P} \left( \sup_{x \in [0,1]} \left| \frac{1}{T} \sum_{t=1}^T W_h(X_{it} - x) \varepsilon_{it}^> \right| > \frac{a_{n,T}}{2} \right) \\ & \leq \sum_{i=1}^n \mathbb{P} \left( |\varepsilon_{it}| > \psi_{n,T} \text{ for some } 1 \leq t \leq T \right) \leq C(nT)^{1-\frac{\theta}{\theta-\delta}} = o(1). \end{aligned}$$

Once more applying Condition 3, it can be seen that

$$\begin{aligned} |\mathbb{E}[W_h(X_{it} - x) \varepsilon_{it}^>]| & \leq \mathbb{E} \left[ W_h(X_{it} - x) \mathbb{E} \left[ \frac{|\varepsilon_{it}|^\theta}{\psi_{n,T}^{\theta-1}} 1(|\varepsilon_{it}| > \psi_{n,T}) \middle| X_{it} \right] \right] \\ & \leq C(nT)^{-\frac{\theta-1}{\theta-\delta}} \end{aligned}$$

with some constant  $C$  independent of  $x$ . Since  $C(nT)^{-\frac{\theta-1}{\theta-\delta}} < a_{n,T}/2$  as the sample size grows large, we arrive at

$$\sum_{i=1}^n \mathbb{P} \left( \sup_{x \in [0,1]} \left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[W_h(X_{it} - x) \varepsilon_{it}^>] \right| > \frac{a_{n,T}}{2} \right) = 0$$

for sufficiently large sample sizes. This yields (S.8).

We next have a closer look at the expression  $\sum_{t=1}^T Z_{it,T}^{\leq}(x)$ . Let  $0 = x_0 < x_1 < \dots < x_L = 1$  be an equidistant grid of points covering the unit interval and set  $L = L_{n,T} = \psi_{n,T}/(a_{n,T}h^2)$ . Exploiting the Lipschitz continuity of the kernel  $W$ , straightforward calculations yield that

$$\max_{1 \leq i \leq n} \sup_{x \in [0,1]} \left| \sum_{t=1}^T Z_{it,T}^{\leq}(x) \right| \leq \max_{1 \leq i \leq n} \max_{1 \leq \ell \leq L} \left| \sum_{t=1}^T Z_{it,T}^{\leq}(x_\ell) \right| + C a_{n,T}. \quad (\text{S.9})$$

We can thus replace the supremum over  $x$  by a maximum over the grid points  $x_\ell$ . Moreover, it holds that

$$\begin{aligned} & \mathbb{P} \left( \max_{1 \leq i \leq n} \max_{1 \leq \ell \leq L} \left| \sum_{t=1}^T Z_{it,T}^{\leq}(x_\ell) \right| > C_0 a_{n,T} \right) \\ & \leq \sum_{i=1}^n \sum_{\ell=1}^L \mathbb{P} \left( \left| \sum_{t=1}^T Z_{it,T}^{\leq}(x_\ell) \right| > C_0 a_{n,T} \right), \end{aligned} \quad (\text{S.10})$$

where  $C_0$  is a sufficiently large constant to be specified later on. In what follows, we show that for each fixed  $x_\ell$ ,

$$\mathbb{P} \left( \left| \sum_{t=1}^T Z_{it,T}^{\leq}(x_\ell) \right| > C_0 a_{n,T} \right) \leq C T^{-r}, \quad (\text{S.11})$$

where the constants  $C$  and  $r$  are independent of  $x_\ell$  and  $r > 0$  can be chosen arbitrarily large provided that  $C_0$  is picked sufficiently large. Plugging (S.11) into (S.10) and combining the result with (S.9), we arrive at

$$\max_{1 \leq i \leq n} \sup_{x \in [0,1]} \left| \sum_{t=1}^T Z_{it,T}^{\leq}(x) \right| = O_p(a_{n,T}), \quad (\text{S.12})$$

which completes the proof.

It thus remains to prove (S.11). To do so, we split the term  $\sum_{t=1}^T Z_{it,T}^{\leq}(x_\ell)$  into blocks as follows:

$$\sum_{t=1}^T Z_{it,T}^{\leq}(x_\ell) = \sum_{s=1}^{q_{n,T}} B_{2s-1} + \sum_{s=1}^{q_{n,T}} B_{2s}$$

with  $B_s = \sum_{t=(s-1)r_{n,T}+1}^{sr_{n,T}} Z_{it,T}^{\leq}(x_\ell)$ , where  $2q_{n,T}$  is the number of blocks and  $r_{n,T} = T/(2q_{n,T})$  is the block length. In particular, we choose the block length such that  $r_{n,T} = O(T^\eta)$  for some small  $\eta > 0$ . With this notation at hand, we get

$$\begin{aligned} \mathbb{P}\left(\left|\sum_{t=1}^T Z_{it,T}^{\leq}(x_\ell)\right| > C_0 a_{n,T}\right) &\leq \mathbb{P}\left(\left|\sum_{s=1}^{q_{n,T}} B_{2s-1}\right| > \frac{C_0}{2} a_{n,T}\right) \\ &\quad + \mathbb{P}\left(\left|\sum_{s=1}^{q_{n,T}} B_{2s}\right| > \frac{C_0}{2} a_{n,T}\right). \end{aligned}$$

As the two terms on the right-hand side can be treated analogously, we focus attention to the first one. By Bradley's lemma (see Lemma 1.2 in Bosq (1998)), we can construct a sequence of random variables  $B_1^*, B_3^*, \dots$  such that (a)  $B_1^*, B_3^*, \dots$  are independent, (b)  $B_{2s-1}$  and  $B_{2s-1}^*$  have the same distribution for each  $s$ , and (c) for  $0 < \mu \leq \|B_{2s-1}\|_\infty$ ,  $\mathbb{P}(|B_{2s-1}^* - B_{2s-1}| > \mu) \leq 18(\|B_{2s-1}\|_\infty/\mu)^{1/2} \alpha(r_{n,T})$ . With these variables, we obtain the bound

$$\mathbb{P}\left(\left|\sum_{s=1}^{q_{n,T}} B_{2s-1}\right| > \frac{C_0}{2} a_{n,T}\right) \leq P_1 + P_2,$$

where

$$\begin{aligned} P_1 &= \mathbb{P}\left(\left|\sum_{s=1}^{q_{n,T}} B_{2s-1}^*\right| > \frac{C_0}{4} a_{n,T}\right) \\ P_2 &= \mathbb{P}\left(\left|\sum_{s=1}^{q_{n,T}} (B_{2s-1} - B_{2s-1}^*)\right| > \frac{C_0}{4} a_{n,T}\right). \end{aligned}$$

Using (c) together with the fact that the mixing coefficients  $\alpha(\cdot)$  decay to zero exponentially fast, it is not difficult to see that  $P_2$  converges to zero at an arbitrarily fast polynomial rate. To deal with  $P_1$ , we make use of the following three facts:

(i) For a random variable  $B$  and  $\lambda > 0$ , Markov's inequality yields that

$$\mathbb{P}(\pm B > \delta) \leq \frac{\mathbb{E} \exp(\pm \lambda B)}{\exp(\lambda \delta)}.$$

(ii) We have that  $|B_{2s-1}| \leq C_B r_{n,T} \psi_{n,T} / (Th)$  for some constant  $C_B > 0$ . Define  $\lambda_{n,T} = Th / (2C_B r_{n,T} \psi_{n,T})$ , which implies that  $\lambda_{n,T} |B_{2s-1}| \leq 1/2$ . As  $\exp(x) \leq 1 + x + x^2$  for  $|x| \leq 1/2$ , we get that

$$\mathbb{E} \left[ \exp(\pm \lambda_{n,T} B_{2s-1}) \right] \leq 1 + \lambda_{n,T}^2 \mathbb{E}[(B_{2s-1})^2] \leq \exp(\lambda_{n,T}^2 \mathbb{E}[(B_{2s-1})^2])$$

along with

$$\mathbb{E} \left[ \exp(\pm \lambda_{n,T} B_{2s-1}^*) \right] \leq \exp(\lambda_{n,T}^2 \mathbb{E}[(B_{2s-1}^*)^2]).$$

(iii) Standard calculations for kernel estimators yield that

$$\sum_{s=1}^{q_{n,T}} \mathbb{E}[(B_{2s-1}^*)^2] \leq \frac{C}{Th}.$$

Using (i)–(iii), we arrive at

$$\begin{aligned} & \mathbb{P} \left( \left| \sum_{s=1}^{q_{n,T}} B_{2s-1}^* \right| > \frac{C_0}{4} a_{n,T} \right) \\ & \leq \mathbb{P} \left( \sum_{s=1}^{q_{n,T}} B_{2s-1}^* > \frac{C_0}{4} a_{n,T} \right) + \mathbb{P} \left( - \sum_{s=1}^{q_{n,T}} B_{2s-1}^* > \frac{C_0}{4} a_{n,T} \right) \\ & \leq \exp \left( - \frac{C_0}{4} \lambda_{n,T} a_{n,T} \right) \left\{ \mathbb{E} \left[ \exp \left( \lambda_{n,T} \sum_{s=1}^{q_{n,T}} B_{2s-1}^* \right) \right] + \mathbb{E} \left[ \exp \left( - \lambda_{n,T} \sum_{s=1}^{q_{n,T}} B_{2s-1}^* \right) \right] \right\} \\ & \leq \exp \left( - \frac{C_0}{4} \lambda_{n,T} a_{n,T} \right) \left\{ \prod_{s=1}^{q_{n,T}} \mathbb{E} \left[ \exp(\lambda_{n,T} B_{2s-1}^*) \right] + \prod_{s=1}^{q_{n,T}} \mathbb{E} \left[ \exp(-\lambda_{n,T} B_{2s-1}^*) \right] \right\} \\ & \leq 2 \exp \left( - \frac{C_0}{4} \lambda_{n,T} a_{n,T} \right) \prod_{s=1}^{q_{n,T}} \exp \left( \lambda_{n,T}^2 \mathbb{E}[(B_{2s-1}^*)^2] \right) \\ & = 2 \exp \left( - \frac{C_0}{4} \lambda_{n,T} a_{n,T} \right) \exp \left( \lambda_{n,T}^2 \sum_{s=1}^{q_{n,T}} \mathbb{E}[(B_{2s-1}^*)^2] \right) \\ & \leq 2 \exp \left( - \frac{C_0}{4} \lambda_{n,T} a_{n,T} + \lambda_{n,T}^2 \frac{C}{Th} \right). \end{aligned}$$

Recalling that  $n/T \leq C$  and  $h \geq cT^{-2/5}$  by assumption, setting  $\theta$  to a value slightly larger than 4 and supposing that  $a_{n,T} = T^{-1/10}$ , it follows that

$$\exp \left( - \frac{C_0}{4} \lambda_{n,T} a_{n,T} + \lambda_{n,T}^2 \frac{C}{Th} \right) \leq T^{-r}, \quad (\text{S.13})$$



where the constant  $r > 0$  can be made arbitrarily large by picking  $C_0$  large enough. If we strengthen Condition 3 to be satisfied for some  $\theta > 20/3$  and choose the block length to be  $r_{n,T} = \sqrt{(Th)/(\psi_{n,T}^2 \log T)}$ , (S.13) even holds for  $a_{n,T} = \sqrt{\log T/(Th)}$ . From (S.13), it immediately follows that  $P_1 \leq CT^{-r}$ , which in turn completes the proof of (S.11).  $\square$

**Proof of (S.5).** The statement follows essentially by the same arguments as those for the proof of (S.4).  $\square$

**Proof of (S.6).** Define  $Z_{it} = (n-1)^{-1} \sum_{j=1, j \neq i}^n (m_j(X_{jt}) + \varepsilon_{jt})$  and write

$$Q_{i,\gamma}(x) = \frac{1}{T} \sum_{t=1}^T W_h(X_{it} - x) Z_{it}. \quad (\text{S.14})$$

By construction, the time series processes  $\{X_{it} : 1 \leq t \leq T\}$  and  $\{Z_{it} : 1 \leq t \leq T\}$  are independent of each other. Moreover, by Theorem 5.2 in Bradley (2005), the process  $\{Z_{it} : 1 \leq t \leq T\}$  is strongly mixing with mixing coefficients that are bounded by  $n\alpha(\ell)$ . (S.6) can thus be shown by applying the arguments from the proof of (S.4) to (S.14).  $\square$

**Proof of (S.7).** The overall strategy is the same as that for the proof of (S.4). There is however one important difference: In the proof of (S.4), we have examined a kernel average of the form  $T^{-1} \sum_{t=1}^T W_h(X_{it} - x) Z_{it}$  with  $Z_{it} = \varepsilon_{it}$ . As the variables  $\varepsilon_{it}$  have unbounded support in general, we have introduced the truncation sequence  $\psi_{n,T}$  and have split  $\varepsilon_{it}$  into the two parts  $\varepsilon_{it}^{\leq}$  and  $\varepsilon_{it}^{>}$ . Here in contrast, we are concerned with the case  $Z_{it} \equiv 1$ . Importantly, the random variables  $Z_{it} \equiv 1$  are bounded, implying that we do not have to truncate them at all. Keeping this in mind and going step by step along the proof of (S.4), we arrive at (S.7).  $\square$

**Proof of Lemma S.2.** Under the conditions of the lemma, it holds that for any pair of indices  $i, j \in G_k$ ,

$$\widehat{\Delta}_{ij} = \int \left( \frac{Q_{i,V}(x) + Q_{i,B}(x)}{\widehat{f}_i(x)} - \frac{Q_{j,V}(x) + Q_{j,B}(x)}{\widehat{f}_j(x)} \right)^2 \pi(x) dx$$

with  $Q_{i,V}(x)$ ,  $Q_{i,B}(x)$  and  $\widehat{f}_i(x)$  defined as in the proof of Lemma S.1. Using the arguments from Lemma S.1, one can show that

$$\begin{aligned} \max_{i \in G_k} \sup_{x \in [0,1]} |Q_{i,V}(x)| &= O_p \left( \sqrt{\frac{\log T}{Th}} \right) \\ \max_{i \in G_k} \sup_{x \in [0,1]} |Q_{i,B}(x) - \mathbb{E}[Q_{i,B}(x)]| &= O_p \left( h \sqrt{\frac{\log T}{Th}} \right) \end{aligned}$$

$$\max_{i \in G_k} \sup_{x \in I_h} |\hat{f}_i(x) - f_i(x)| = O_p\left(\sqrt{\frac{\log T}{Th}} + h^2\right)$$

and

$$\mathbb{E}[Q_{i,B}(x)] = h^2 \left( \int W(\varphi) \varphi^2 d\varphi \right) \left( m'_i(x) f'_i(x) + \frac{m''_i(x) f_i(x)}{2} \right) + O(h^3)$$

uniformly for  $i \in G_k$  and  $x \in I_h$ . Applying these uniform convergence results and noting that  $\max_{i,j \in G_k} \hat{\Delta}_{ij} = \max_{i,j \in G_k, i < j} \hat{\Delta}_{ij}$ , it is not difficult to see that

$$\max_{i,j \in G_k} \hat{\Delta}_{ij} = \max_{\substack{i,j \in G_k \\ i < j}} \int \left( \frac{Q_{i,V}(x)}{f_i(x)} - \frac{Q_{j,V}(x)}{f_j(x)} \right)^2 \pi(x) dx + o_p\left(\frac{1}{Th^{1/2}}\right).$$

Next define

$$U_{i,T} = \sum_{s,t=1}^T a_{st}^{(i)} \varepsilon_{is} \varepsilon_{it},$$

where  $a_{st}^{(i)} = T^{-2} \int W_h(X_{is} - x) W_h(X_{it} - x) \pi(x) / f_i^2(x) dx$  for  $s \neq t$  and  $a_{st}^{(i)} = 0$  for  $s = t$ . Similarly, for  $i \neq j$ , let

$$U_{ij,T} = \sum_{s,t=1}^T a_{st}^{(ij)} \varepsilon_{is} \varepsilon_{jt}$$

with  $a_{st}^{(ij)} = T^{-2} \int W_h(X_{is} - x) W_h(X_{jt} - x) \pi(x) / (f_i(x) f_j(x)) dx$  and define

$$B_{i,T} = \int \frac{1}{T^2} \sum_{t=1}^T W_h^2(X_{it} - x) \varepsilon_{it}^2 \frac{\pi(x)}{f_i^2(x)} dx.$$

With these definitions at hand, we can write

$$\max_{\substack{i,j \in G_k \\ i < j}} \int \left( \frac{Q_{i,V}(x)}{f_i(x)} - \frac{Q_{j,V}(x)}{f_j(x)} \right)^2 \pi(x) dx = \max_{\substack{i,j \in G_k \\ i < j}} \{U_{i,T} - 2U_{ij,T} + U_{j,T} + B_{i,T} + B_{j,T}\}.$$

Below we show that

$$\max_{1 \leq i \leq n} |U_{i,T}| = O_p\left(\frac{\log T}{Th^{1/2}}\right) \quad (\text{S.15})$$

$$\max_{1 \leq i < j \leq n} |U_{ij,T}| = O_p\left(\frac{\log T}{Th^{1/2}}\right). \quad (\text{S.16})$$

Moreover, similar arguments as those for the proof of Lemma S.1 yield that

$$\max_{1 \leq i \leq n} |B_{i,T} - \mathbb{E}[B_{i,T}]| = o_p\left(\frac{1}{Th^{1/2}}\right) \quad (\text{S.17})$$

$$\mathbb{E}[B_{i,T}] = \frac{1}{Th} \left( \int W^2(\varphi) d\varphi \right) \int \frac{\sigma_i^2(x) \pi(x)}{f_i(x)} dx + O\left(\frac{1}{T}\right) \quad (\text{S.18})$$

uniformly in  $i$ . Combining (S.15)–(S.18) and noting that  $\mathbb{E}[B_{i,T}] + \mathbb{E}[B_{j,T}] = \mathcal{B}_{ij}/(Th) + O(T^{-1})$  uniformly in  $i$  and  $j$ , we arrive at

$$\begin{aligned} \max_{i,j \in G_k} \widehat{\Delta}_{ij} &= \max_{\substack{i,j \in G_k \\ i < j}} \int \left( \frac{Q_{i,V}(x)}{f_i(x)} - \frac{Q_{j,V}(x)}{f_j(x)} \right)^2 \pi(x) dx + o_p\left(\frac{1}{Th^{1/2}}\right) \\ &= \max_{\substack{i,j \in G_k \\ i < j}} \frac{\mathcal{B}_{ij}}{Th} + O_p\left(\frac{\log T}{Th^{1/2}}\right). \end{aligned}$$

To complete the proof of Lemma S.2, it thus remains to verify (S.15) and (S.16).  $\square$

**Proof of (S.15).** Define the matrix  $A_T^{(i)} = (|a_{st}^{(i)}|)_{s,t=1}^T$  and  $\Lambda_T^{(i)} = \sum_{s,t=1}^T (a_{st}^{(i)})^2$ . We first show that

$$\max_{1 \leq i \leq n} \|A_T^{(i)}\| = O_p\left(\frac{1}{T}\right) \quad (\text{S.19})$$

$$\max_{1 \leq i \leq n} \Lambda_T^{(i)} = O_p\left(\frac{1}{T^2 h}\right), \quad (\text{S.20})$$

where  $\|A_T^{(i)}\|$  denotes the spectral norm of  $A_T^{(i)}$ . By definition,  $\|A_T^{(i)}\|$  is the largest absolute eigenvalue of  $A_T^{(i)}$ . As the diagonal elements  $|a_{tt}^{(i)}|$  of  $A_T^{(i)}$  are all zero, Gerschgorin's theorem says that the largest absolute eigenvalue of  $A_T^{(i)}$  is bounded by

$$\bar{\lambda}_T^{(i)} = \max_{1 \leq s \leq T} \sum_{t=1}^T |a_{st}^{(i)}|.$$

Standard calculations yield that

$$\bar{\lambda}_T^{(i)} \leq \frac{C}{T} \max_{1 \leq s \leq T} \frac{1}{T} \sum_{t=1}^T \mathcal{W}_h(X_{is} - X_{it}),$$

where  $\mathcal{W}_h(x) = h^{-1}\mathcal{W}(x/h)$  and  $\mathcal{W}(x) = \int_{-C_1}^{C_1} W(x + \varphi) d\varphi$ . One can easily show that (a)  $|\mathcal{W}(x)| \leq C$ , (b)  $\mathcal{W}(x) = 0$  for all  $|x| > 2C_1$ , and (c)  $|\mathcal{W}(x) - \mathcal{W}(x')| \leq L|x - x'|$  for some constant  $L$ . Hence, similar arguments as those from Lemma S.1 yield that

$$\begin{aligned} \bar{\lambda}_T^{(i)} &\leq \frac{C}{T} \sup_{x \in [0,1]} \frac{1}{T} \sum_{t=1}^T \mathcal{W}_h(x - X_{it}) \\ &\leq \frac{C}{T} \sup_{x \in [0,1]} \left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\mathcal{W}_h(x - X_{it})] \right| \\ &\quad + \frac{C}{T} \sup_{x \in [0,1]} \left| \frac{1}{T} \sum_{t=1}^T \{ \mathcal{W}_h(x - X_{it}) - \mathbb{E}[\mathcal{W}_h(x - X_{it})] \} \right| \\ &= O\left(\frac{1}{T}\right) + O_p\left(\frac{1}{T} \sqrt{\frac{\log T}{Th}}\right) \end{aligned}$$

uniformly over  $i$ . As a result, we get that  $\max_{1 \leq i \leq n} \|A_T^{(i)}\| \leq \max_{1 \leq i \leq n} \bar{\lambda}_T^{(i)} = O_p(T^{-1})$ , thus completing the proof of (S.19). To see (S.20), note that  $\int W_h(X_{is} - x)W_h(X_{it} - x)\pi(x)/f_i^2(x)dx \leq C/h$ . Keeping this in mind, we obtain that

$$\begin{aligned}\Lambda_T^{(i)} &= \frac{1}{T^4} \sum_{\substack{s,t=1 \\ s \neq t}}^T \left\{ \int W_h(X_{is} - x)W_h(X_{it} - x) \frac{\pi(x)}{f_i^2(x)} dx \right\}^2 \\ &\leq \frac{C}{T^4 h} \sum_{s,t=1}^T \int W_h(X_{is} - x)W_h(X_{it} - x) \frac{\pi(x)}{f_i^2(x)} dx \\ &= \frac{C}{T^2 h} \int \left( \frac{1}{T} \sum_{s=1}^T W_h(X_{is} - x) \right) \left( \frac{1}{T} \sum_{t=1}^T W_h(X_{it} - x) \right) \frac{\pi(x)}{f_i^2(x)} dx \\ &= O_p\left(\frac{1}{T^2 h}\right)\end{aligned}$$

uniformly in  $i$ , taking into account that  $T^{-1} \sum_{t=1}^T W_h(X_{it} - x) = O_p(1)$  uniformly over  $i$  and  $x$ .

We now let  $\mathcal{X}_{n,T} = (X_{11}, \dots, X_{1T}, X_{21}, \dots, X_{2T}, \dots, X_{n1}, \dots, X_{nT})$  be the vector of the regressors  $X_{it}$  and define the event

$$E_T = \left\{ \mathcal{X}_{n,T} : \max_{1 \leq i \leq n} \|A_T^{(i)}\| \leq \frac{\log T}{T} \text{ and } \max_{1 \leq i \leq n} \Lambda_T^{(i)} \leq \frac{\log T}{T^2 h} \right\}.$$

By (S.19) and (S.20), it holds that  $\mathbb{P}(E_T) \rightarrow 1$ . Hence,

$$\begin{aligned}\mathbb{P}\left(\max_{1 \leq i \leq n} |U_{i,T}| > C_U \frac{\log T}{Th^{1/2}}\right) &= \mathbb{P}\left(\max_{1 \leq i \leq n} |U_{i,T}| > C_U \frac{\log T}{Th^{1/2}}, E_T\right) + o(1) \\ &\leq \sum_{i=1}^n \mathbb{P}\left(|U_{i,T}| > C_U \frac{\log T}{Th^{1/2}}, E_T\right) + o(1) \\ &= \sum_{i=1}^n \mathbb{P}\left(1(E_T) |U_{i,T}| > C_U \frac{\log T}{Th^{1/2}}\right) + o(1).\end{aligned}$$

We further write

$$\mathbb{P}\left(1(E_T) |U_{i,T}| > C_U \frac{\log T}{Th^{1/2}}\right) = \mathbb{E}\left[\mathbb{P}\left(1(E_T) |U_{i,T}| > C_U \frac{\log T}{Th^{1/2}} \mid \mathcal{X}_{n,T}\right)\right]$$

and derive an exponential bound on the conditional probability  $\mathbb{P}(1(E_T) |U_{i,T}| > C_U \log T / (Th^{1/2}) | \mathcal{X}_{n,T})$ . To do so, we make use of the following result, which is immediately implied by the proof of the theorem in Wright (1973).

**Theorem.** *Define*

$$U_T = \sum_{s,t=-T}^T a_{st}(\eta_s \eta_t - \mathbb{E}[\eta_s \eta_t])$$

*and suppose that the following conditions are satisfied:*

(i)  $\{\eta_t : -T \leq t \leq T\}$  is a sequence of independent random variables with zero means. For some constants  $M, \gamma > 0$ ,  $\mathbb{P}(|\eta_t| \geq c) \leq M \int_c^\infty \exp(-\gamma r^2) dr$  for all  $-T \leq t \leq T$  and all  $c \geq 0$ .

(ii) For  $-T \leq s, t \leq T$ ,  $a_{st}$  are real numbers with  $a_{st} = a_{ts}$  and  $\Lambda_T = \sum_{s,t=-T}^T a_{st}^2 \leq C < \infty$ . Let  $A_T = (|a_{st}|)_{s,t=-T}^T$  and denote the spectral norm of  $A_T$  by  $\|A_T\|$ .

Then there exist positive constants  $C_a$  and  $C_b$  depending only on  $M$  and  $\gamma$  such that for every  $\delta > 0$ ,

$$\mathbb{P}(U_T > \delta) \leq \exp \left( - \min \left\{ \frac{C_a \delta}{\|A_T\|}, \frac{C_b \delta^2}{\Lambda_T} \right\} \right).$$

Setting  $a_{st}^{(i)} = 0$  whenever  $s < 1$  or  $t < 1$ , we can write  $U_{i,T} = \sum_{s,t=-T}^T a_{st}^{(i)} \varepsilon_{is} \varepsilon_{it}$  and directly apply the above theorem. This yields

$$\begin{aligned} & \mathbb{P} \left( 1(E_T) |U_{i,T}| > C_U \frac{\log T}{Th^{1/2}} \mid \mathcal{X}_{n,T} \right) \\ & \leq \exp \left( - \min \left\{ \frac{C_a C_U \log T / (Th^{1/2})}{\log T / T}, \frac{C_b C_U^2 (\log T / (Th^{1/2}))^2}{\log T / (T^2 h)} \right\} \right) \\ & = \exp \left( - C_b C_U^2 \log T \right) = T^{-C_b C_U^2} \end{aligned}$$

for sufficiently large sample sizes  $T$ . As a result,

$$\mathbb{P} \left( \max_{1 \leq i \leq n} |U_{i,T}| > C_U \frac{\log T}{Th^{1/2}} \right) \leq n T^{-C_b C_U^2} + o(1) = o(1)$$

for  $C_U$  chosen sufficiently large. □

**Proof of (S.16).** First of all, note that we can write

$$U_{ij,T} = \int \frac{Q_{i,V}(x) Q_{j,V}(x)}{f_i(x) f_j(x)} \pi(x) dx$$

with  $Q_{i,V}(x) = T^{-1} \sum_{t=1}^T W_h(X_{it} - x) \varepsilon_{it}$ . The arguments from the proof of Lemma S.1 show that

$$\mathbb{P} \left( \max_{1 \leq i \leq n} \sup_{x \in [0,1]} |Q_{i,V}(x)| > C_Q \sqrt{\frac{\log T}{Th}} \right) = o(1)$$

for  $C_Q$  chosen sufficiently large. Now let  $E_T$  be the event that  $\max_i \sup_x |Q_{i,V}(x)| \leq C_Q \sqrt{\log T / (Th)}$  and  $E_{j,T}$  the event that  $\sup_x |Q_{j,V}(x)| \leq C_Q \sqrt{\log T / (Th)}$ . Then

$$\begin{aligned} \mathbb{P} \left( \max_{1 \leq i < j \leq n} |U_{ij,T}| > C_U \frac{\log T}{Th^{1/2}} \right) &= \mathbb{P} \left( \max_{1 \leq i < j \leq n} |U_{ij,T}| > C_U \frac{\log T}{Th^{1/2}}, E_T \right) + o(1) \\ &\leq \sum_{1 \leq i < j \leq n} \mathbb{P} \left( |U_{ij,T}| > C_U \frac{\log T}{Th^{1/2}}, E_T \right) + o(1) \end{aligned}$$

and

$$\begin{aligned}
\mathbb{P}\left(|U_{ij,T}| > C_U \frac{\log T}{Th^{1/2}}, E_T\right) &= \mathbb{P}\left(1(E_T)|U_{ij,T}| > C_U \frac{\log T}{Th^{1/2}}\right) \\
&\leq \mathbb{P}\left(1(E_{j,T})|U_{ij,T}| > C_U \frac{\log T}{Th^{1/2}}\right) \\
&= \mathbb{P}\left(\left|\frac{1}{T} \sum_{t=1}^T w_{ij,T} \varepsilon_{it}\right| > C_U \frac{\log T}{Th^{1/2}}\right),
\end{aligned}$$

where we set

$$w_{ij,T} = \int \frac{W_h(X_{it} - x)}{f_i(x)} \frac{Q_{j,V}(x)1(E_{j,T})}{f_j(x)} \pi(x) dx.$$

Noting that  $w_{ij,T} \leq C \sqrt{\log T / (Th)}$ , one can show that

$$\mathbb{P}\left(\left|\frac{1}{T} \sum_{t=1}^T w_{ij,T} \varepsilon_{it}\right| > C_U \frac{\log T}{Th^{1/2}}\right) \leq CT^{-r},$$

where  $r > 0$  can be made arbitrarily large by choosing  $C_U$  large enough. This implies that

$$\mathbb{P}\left(\max_{1 \leq i < j \leq n} |U_{ij,T}| > C_U \frac{\log T}{Th^{1/2}}\right) = o(1)$$

for  $C_U$  sufficiently large.  $\square$

## Proof of Theorem 1

Let  $S \subseteq \{1, \dots, n\}$  be some index set with  $n_S = |S|$ , pick an index  $i \in S$ , and let  $G \subseteq S$  be the class to which  $i$  belongs. As seen in Subsection 2.2 of the paper, the group  $G$  has the form  $G = \{(1), \dots, (p)\}$ , where  $\Delta_{i(1)} = \dots = \Delta_{i(p)} < \Delta_{i(p+1)} \leq \dots \leq \Delta_{i(n_S)}$  are the ordered  $L_2$ -distances. Denoting the ordered estimated distances by  $\hat{\Delta}_{i[1]} \leq \hat{\Delta}_{i[2]} \leq \dots \leq \hat{\Delta}_{i[n_S]}$ , we estimate  $G$  by  $\hat{G} = \{[1], \dots, [\hat{p}]\}$  with  $\hat{p}$  defined in (2.3). In what follows, we show that

$$\mathbb{P}\left(\{[1], \dots, [\hat{p}]\} \neq \{(1), \dots, (p)\}\right) = o(1). \quad (\text{S.21})$$

With the help of (S.21) and some straightforward additional arguments, the statements of Theorem 1 can be easily inferred. For the proof of (S.21), it suffices to show that

$$\mathbb{P}\left(\{[1], \dots, [p]\} \neq \{(1), \dots, (p)\}\right) = o(1) \quad (\text{S.22})$$

$$\mathbb{P}(\hat{p} \neq p) = o(1). \quad (\text{S.23})$$

These two statements can be verified as follows: By (S.2), it holds that  $\hat{\Delta}_{i(j)} - \Delta_{i(j)} = o_p(1)$  uniformly over  $j$ . As  $\Delta_{i(j)} = 0$  for all  $j \leq p$  and  $\Delta_{i(j)} \geq c$  for all  $j > p$  and some

constant  $c > 0$ , we obtain that

$$\max_{j \leq p} \widehat{\Delta}_{i(j)} = o_p(1) \quad \text{and} \quad \min_{j > p} \widehat{\Delta}_{i(j)} \geq c + o_p(1). \quad (\text{S.24})$$

This immediately implies that the ordered estimates  $\widehat{\Delta}_{i[j]}$  have the same property, i.e.,

$$\max_{j \leq p} \widehat{\Delta}_{i[j]} = o_p(1) \quad \text{and} \quad \min_{j > p} \widehat{\Delta}_{i[j]} \geq c + o_p(1). \quad (\text{S.25})$$

From (S.24) and (S.25), it is obvious that the index sets  $\{[1], \dots, [p]\}$  and  $\{(1), \dots, (p)\}$  coincide with probability tending to one, which is the statement of (S.22). From (S.22), it follows that  $\max_{j \in G} \widehat{\Delta}_{ij} = \widehat{\Delta}_{i[p]}$  with probability tending to one. Moreover, as the threshold parameter  $\tau_{n,T}$  satisfies Condition 6,  $\widehat{\Delta}_{i[p]} \leq \tau_{n,T}$  with probability approaching one. Finally, by (S.25),  $\widehat{\Delta}_{i[p+1]} > \tau_{n,T}$  with probability approaching one as well. We thus arrive at

$$\mathbb{P}(\widehat{\Delta}_{i[p]} \leq \tau_{n,T} \text{ and } \widehat{\Delta}_{i[p+1]} > \tau_{n,T}) \rightarrow 1,$$

which immediately implies that  $\mathbb{P}(\widehat{p} = p) \rightarrow 1$ .  $\square$

## Proof of Theorem 2

As  $\widehat{K} = K$  with probability tending to one, we can neglect the estimation error in  $\widehat{K}$  and treat  $K$  as known. With the help of Lemma S.1, it is straightforward to see that

$$\int (\widehat{m}_i(x) - \widehat{g}_k^{[1]}(x))^2 \pi(x) dx = \int (m_i(x) - g_k(x))^2 \pi(x) dx + o_p(1)$$

uniformly over  $i$  and  $k$ , or put differently,

$$\max_{1 \leq k \leq K} \max_{1 \leq i \leq n} |\Delta(\widehat{m}_i, \widehat{g}_k^{[1]}) - \Delta(m_i, g_k)| = o_p(1). \quad (\text{S.26})$$

By construction, the index  $i$  is assigned to the group  $G_k^{[1]}$  in the first step of the  $k$ -means algorithm if  $\widehat{d}_k(i) = \Delta(\widehat{m}_i, \widehat{g}_k^{[1]})$  is minimal, i.e., if  $\widehat{d}_k(i) = \min_{1 \leq k' \leq K} \widehat{d}_{k'}(i)$ . By (S.26), we know that

$$\widehat{d}_k(i) = \begin{cases} \widehat{r}_k(i) & \text{if } i \in G_k \\ \Delta(m_i, g_k) + \widehat{r}_k(i) & \text{if } i \notin G_k, \end{cases} \quad (\text{S.27})$$

where the remainder term  $\widehat{r}_k(i)$  has the property that  $\max_{1 \leq k \leq K} \max_{1 \leq i \leq n} |\widehat{r}_k(i)| = o_p(1)$ . Since  $\min_{1 \leq k \leq K} \min_{i \notin G_k} \Delta(m_i, g_k) \geq \Delta_{\min} > 0$  for some positive constant  $\Delta_{\min}$ , (S.27) implies that

$$\mathbb{P}(\{G_k^{[1]} : 1 \leq k \leq K\} \neq \{G_k : 1 \leq k \leq K\}) = o(1).$$

Hence, with probability tending to one, our  $k$ -means clustering algorithm converges already after the first iteration step and produces estimates which coincide with the classes  $G_k$  for  $1 \leq k \leq K$ .  $\square$

### Proof of Theorem 3

We focus attention on the proof of the distribution result (3.2). The convergence result (3.1) follows by slightly modifying the arguments of the proof. In a first step, we replace the estimator  $\widehat{g}_k$  by the infeasible version

$$\widehat{g}_k^*(x) = \frac{1}{n_k} \sum_{i \in G_k} \widehat{m}_i(x)$$

and show that the difference between the two estimators is asymptotically negligible: For any null sequence  $\{a_{n,T}\}$  of positive numbers, it holds that

$$\begin{aligned} & \mathbb{P}\left(\left|\widehat{g}_k(x) - \widehat{g}_k^*(x)\right| > a_{n,T}\right) \\ & \leq \mathbb{P}\left(\left|\widehat{g}_k(x) - \widehat{g}_k^*(x)\right| > a_{n,T}, \widehat{G}_k = G_k\right) + \mathbb{P}(\widehat{G}_k \neq G_k) = o(1), \end{aligned}$$

since the first probability on the right-hand side is equal to zero by definition of  $\widehat{g}_k$  and  $\widehat{g}_k^*$  and the second one is of the order  $o(1)$  by Theorem 1. Hence,  $|\widehat{g}_k(x) - \widehat{g}_k^*(x)| = O_p(a_{n,T})$  for any null sequence  $\{a_{n,T}\}$  of positive numbers, which in turn implies that

$$\sqrt{\widehat{n}_k Th}(\widehat{g}_k(x) - g_k(x)) = \sqrt{\widehat{n}_k Th}(\widehat{g}_k^*(x) - g_k(x)) + o_p(1).$$

The difference between  $\widehat{g}_k$  and  $\widehat{g}_k^*$  can thus be asymptotically ignored.

To complete the proof of Theorem 3, we derive the limit distribution of the term  $\sqrt{\widehat{n}_k Th}(\widehat{g}_k^*(x) - g_k(x))$ : Since  $\mathbb{P}(\widehat{n}_k \neq n_k) = o(1)$  by Theorem 1, it holds that  $\sqrt{\widehat{n}_k Th}(\widehat{g}_k^*(x) - g_k(x)) = \sqrt{n_k Th}(\widehat{g}_k^*(x) - g_k(x)) + o_p(1)$ . It thus suffices to compute the limit distribution of  $\sqrt{n_k Th}(\widehat{g}_k^*(x) - g_k(x))$ . To do so, write

$$\widehat{m}_i(x) - m_i(x) = [Q_{i,V}(x) + Q_{i,B}(x) - Q_{i,\gamma}(x)] / \widehat{f}_i(x) - \overline{Q}_i + \overline{\overline{Q}}_i,$$

where  $Q_{i,V}(x)$ ,  $Q_{i,B}(x)$ ,  $Q_{i,\gamma}(x)$  along with  $\overline{Q}_i$ ,  $\overline{\overline{Q}}_i$  and  $\widehat{f}_i(x)$  are defined in the proof of Lemma S.1. With this notation at hand, we obtain that

$$\begin{aligned} & \sqrt{n_k Th}(\widehat{g}_k^*(x) - g_k(x)) \\ & = \sqrt{n_k Th} \left\{ \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,V}(x)}{\widehat{f}_i(x)} + \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,B}(x)}{\widehat{f}_i(x)} - \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\gamma}(x)}{\widehat{f}_i(x)} - \frac{1}{n_k} \sum_{i \in G_k} (\overline{Q}_i - \overline{\overline{Q}}_i) \right\} \\ & = \sqrt{n_k Th} \left\{ \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,V}(x)}{\widehat{f}_i(x)} + \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,B}(x)}{\widehat{f}_i(x)} - \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\gamma}(x)}{\widehat{f}_i(x)} \right\} + o_p(1), \end{aligned}$$



the last line following by standard calculations. In the sequel, we show that

$$\frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\gamma}(x)}{\widehat{f}_i(x)} = o_p\left(\frac{1}{\sqrt{n_k T h}}\right) \quad (\text{S.28})$$

$$\frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,V}(x)}{\widehat{f}_i(x)} = \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,V}(x)}{f_i(x)} + o_p\left(\frac{1}{\sqrt{n_k T h}}\right) \quad (\text{S.29})$$

$$\frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,B}(x)}{\widehat{f}_i(x)} = \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,B}(x)}{f_i(x)} + o_p\left(\frac{1}{\sqrt{n_k T h}}\right). \quad (\text{S.30})$$

(S.28)–(S.30) allow us to conclude that

$$\begin{aligned} & \sqrt{n_k T h} (\widehat{g}_k^*(x) - g_k(x)) \\ &= \sqrt{n_k T h} \left\{ \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,V}(x)}{f_i(x)} + \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,B}(x)}{f_i(x)} \right\} + o_p(1) \\ &= \sqrt{n_k T h} \left( \frac{1}{n_k T} \sum_{i \in G_k} \sum_{t=1}^T \frac{W_h(X_{it} - x)}{f_i(x)} \varepsilon_{it} \right) \\ & \quad + \sqrt{n_k T h} \left( \frac{1}{n_k T} \sum_{i \in G_k} \sum_{t=1}^T \frac{W_h(X_{it} - x)}{f_i(x)} [m_i(X_{it}) - m_i(x)] \right) + o_p(1). \end{aligned}$$

With the help of a standard central limit theorem, the first term on the right-hand side can be shown to weakly converge to a normal distribution with mean zero and variance  $V_k(x)$ . Moreover, standard bias calculations yield that the second term converges in probability to the bias expression  $B_k(x)$ . This completes the proof.  $\square$

**Proof of (S.28).** In a first step, we show that the term

$$R_\gamma := \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\gamma}(x)}{\widehat{f}_i(x)} - \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\gamma}(x)}{\mathbb{E}[\widehat{f}_i(x)]}$$

is of the order

$$R_\gamma = o_p\left(\frac{1}{\sqrt{n_k T h}}\right). \quad (\text{S.31})$$

To do so, we write  $R_\gamma = R_{\gamma,1} + R_{\gamma,2}$ , where

$$\begin{aligned} R_{\gamma,1} &= \frac{1}{n_k} \sum_{i \in G_k} \frac{\mathbb{E}[\widehat{f}_i(x)] - \widehat{f}_i(x)}{\mathbb{E}[\widehat{f}_i(x)]^2} Q_{i,\gamma}(x) \\ R_{\gamma,2} &= \frac{1}{n_k} \sum_{i \in G_k} \frac{(\mathbb{E}[\widehat{f}_i(x)] - \widehat{f}_i(x))^2}{\mathbb{E}[\widehat{f}_i(x)]^2 \widehat{f}_i(x)} Q_{i,\gamma}(x). \end{aligned}$$

Defining  $Z_{it}(x) = \mathbb{E}[W_h(X_{it} - x)] - W_h(X_{it} - x)$ , the term  $R_{\gamma,1}$  can be expressed as

$$R_{\gamma,1} = \frac{1}{n_k} \sum_{i \in G_k} \frac{1}{\mathbb{E}[\widehat{f}_i(x)]^2} \left\{ \frac{1}{T} \sum_{t=1}^T Z_{it}(x) \right\} \\ \times \left\{ \frac{1}{T} \sum_{t=1}^T W_h(X_{it} - x) \left( \frac{1}{n-1} \sum_{\substack{j=1 \\ j \neq i}}^n [m_j(X_{jt}) + \varepsilon_{jt}] \right) \right\}.$$

We thus obtain that

$$\mathbb{E}[R_{\gamma,1}^2] = \frac{1}{n_k^2(n-1)^2} \sum_{i,i' \in G_k} \sum_{\substack{j \neq i \\ j' \neq i'}} \frac{1}{\mathbb{E}[\widehat{f}_i(x)]^2} \frac{1}{\mathbb{E}[\widehat{f}_{i'}(x)]^2} \\ \times \left( \frac{1}{T^4} \sum_{t,t',s,s'=1}^T \Psi_{i,i',j,j',t,t',s,s'}(x) \right), \quad (\text{S.32})$$

where we use the shorthand

$$\Psi_{i,i',j,j',t,t',s,s'}(x) = \mathbb{E}[Z_{it}(x)W_h(X_{is} - x)\{m_j(X_{js}) + \varepsilon_{js}\}] \\ \times Z_{i't'}(x)W_h(X_{i's'} - x)\{m_{j'}(X_{j's'}) + \varepsilon_{j's'}\}].$$

Importantly, the expressions  $\Psi_{i,i',j,j',t,t',s,s'}(x)$  in (S.32) have the following property:  $\Psi_{i,i',j,j',t,t',s,s'}(x) \neq 0$  only if (a)  $i = j'$  and  $i' = j$  or (b)  $j = j'$ . Exploiting the mixing assumptions of Condition 1 by means of Davydov's inequality (see Corollary 1.1 in Bosq (1998)), we can show that in case (a),  $|T^{-4} \sum_{t,t',s,s'=1}^T \psi_{i,i',j,j',t,t',s,s'}(x)| \leq C(\log T)^2/(Th)^2$  and in case (b),

$$\left| \frac{1}{T^4} \sum_{t,t',s,s'=1}^T \psi_{i,i',j,j',t,t',s,s'}(x) \right| \leq \begin{cases} C(\log T)^3/(T^3h^2) & \text{for } i \neq i' \\ C(\log T)^2/(T^2h^3) & \text{for } i = i'. \end{cases}$$

Plugging these bounds into (S.32), we immediately arrive at  $R_{\gamma,1} = o_p(1/\sqrt{n_k Th})$ . Furthermore, with the help of Hölder's inequality and (S.7), we obtain that

$$R_{\gamma,2} \leq \left\{ \max_{1 \leq i \leq n} \frac{(\mathbb{E}[\widehat{f}_i(x)] - \widehat{f}_i(x))^2}{\mathbb{E}[\widehat{f}_i(x)]^2 \widehat{f}_i(x)} \right\} \left\{ \frac{1}{n_k} \sum_{i \in G_k} \left( \frac{1}{T} \sum_{t=1}^T W_h^{4/3}(X_{it} - x) \right)^{3/4} \right. \\ \left. \times \left( \frac{1}{T} \sum_{t=1}^T \left( \frac{1}{n-1} \sum_{\substack{j=1 \\ j \neq i}}^n [m_j(X_{jt}) + \varepsilon_{jt}] \right)^4 \right)^{1/4} \right\} \\ = O_p \left( \left( \sqrt{\frac{\log T}{Th}} \right)^2 \frac{1}{h^{1/4}(n-1)^{1/2}} \right) = o_p \left( \frac{1}{\sqrt{n_k Th}} \right),$$

which completes the proof of (S.31).

In the next step, we show that

$$\frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\gamma}(x)}{\mathbb{E}[\widehat{f}_i(x)]} = o_p\left(\frac{1}{\sqrt{n_k Th}}\right). \quad (\text{S.33})$$

To do so, we derive the convergence rate of the second moment

$$\begin{aligned} \mathbb{E} \left[ \left\{ \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\gamma}(x)}{\mathbb{E}[\widehat{f}_i(x)]} \right\}^2 \right] &= \frac{1}{n_k^2(n-1)^2} \sum_{i,i' \in G_k} \sum_{\substack{j \neq i \\ j' \neq i'}} \frac{1}{\mathbb{E}[\widehat{f}_i(x)]} \frac{1}{\mathbb{E}[\widehat{f}_{i'}(x)]} \\ &\quad \times \left( \frac{1}{T^2} \sum_{t,t'=1}^T \Psi_{i,i',j,j',t,t'}(x) \right), \end{aligned} \quad (\text{S.34})$$

where  $\Psi_{i,i',j,j',t,t'}(x) = \mathbb{E}[W_h(X_{it} - x)\{m_j(X_{jt}) + \varepsilon_{jt}\}W_h(X_{i't'} - x)\{m_{j'}(X_{j't'}) + \varepsilon_{j't'}\}]$ . Similarly as above,  $\Psi_{i,i',j,j',t,t'}(x) \neq 0$  only if (a)  $i = j'$  and  $i' = j$  or (b)  $j = j'$ . Applying Davydov's inequality once again, we get that in case (a),  $|T^{-2} \sum_{t,t'=1}^T \Psi_{i,i',j,j',t,t'}(x)| \leq C \log T/T$  and in case (b),

$$\left| \frac{1}{T^2} \sum_{t,t'=1}^T \Psi_{i,i',j,j',t,t'}(x) \right| \leq \begin{cases} C/T & \text{for } i \neq i' \\ C/(Th) & \text{for } i = i'. \end{cases}$$

Plugging these bounds into (S.34), we easily arrive at (S.33). The statement (S.28) now follows upon combining (S.31) with (S.33).  $\square$

**Proof of (S.29) and (S.30).** By arguments similar to those for (S.28),

$$\frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\ell}(x)}{\widehat{f}_i(x)} - \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\ell}(x)}{\mathbb{E}[\widehat{f}_i(x)]} = o_p\left(\frac{1}{\sqrt{n_k Th}}\right) \quad (\text{S.35})$$

for  $\ell \in \{V, B\}$ . With the help of standard bias calculations, we further obtain that

$$\frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\ell}(x)}{\mathbb{E}[\widehat{f}_i(x)]} - \frac{1}{n_k} \sum_{i \in G_k} \frac{Q_{i,\ell}(x)}{f_i(x)} = o_p\left(\frac{1}{\sqrt{n_k Th}}\right). \quad (\text{S.36})$$

Combining (S.35) and (S.36) completes the proof.  $\square$

## References

- BOSQ, D. (1998). *Nonparametric statistics for stochastic processes*. New York, Springer.
- BRADLEY, R. C. (2005). Basic properties of strong mixing conditions. A survey and some open questions. *Probability Surveys*, **2** 107–144.
- HANSEN, B. (2008). Uniform convergence rates for kernel estimation with dependent data. *Econometric Theory*, **24** 726–748.

- HÄRDLE, W., HALL, P. and MARRON, J. S. (1988). How far are automatically chosen regression smoothing parameters from their optimum? *J. Am. Statist. Ass.*, **83** 86–95.
- MASRY, E. (1996). Multivariate local polynomial regression for time series: uniform strong consistency and rates. *J. Time Ser. Anal.*, **17** 571–599.
- WRIGHT, F. T. (1973). A bound on tail probabilities for quadratic forms in independent random variables whose distributions are not necessarily symmetric. *Ann. Prob.*, **1** 1068–1070.