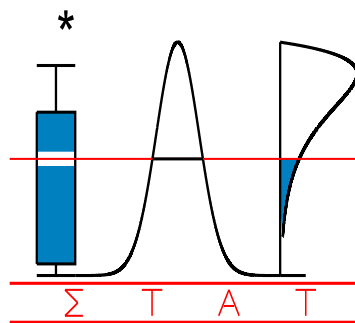


T E C H N I C A L
R E P O R T

0508

**BACKFITTING VERSUS PROFILING
IN GENERAL CRITERION FUNCTIONS**

VAN KEILEGOM, I. and R. J. CARROLL



I A P S T A T I S T I C S
N E T W O R K

INTERUNIVERSITY ATTRACTION POLE

<http://www.stat.ucl.ac.be/IAP>

Backfitting versus Profiling in General Criterion Functions

Ingrid Van Keilegom¹
Institut de Statistique
Université catholique de Louvain
Voie du Roman Pays 20
B-1348 Louvain-la-Neuve, Belgium
vankeilegom@stat.ucl.ac.be

Raymond J. Carroll²
Department of Statistics
Texas A&M University
College Station TX 77843-3143 USA
carroll@stat.tamu.edu

Abstract

We study the backfitting and profile methods for general criterion functions that depend on a parameter of interest β and a nuisance function θ . We show that when different amounts of smoothing are employed for each method to estimate the function θ , the two estimation procedures produce asymptotically the same estimator of β , even when the criterion functions are non-smooth in β and/or θ . The results are applied to a partial linear median regression model and a change point model, both examples of non-smooth criterion functions.

Key Words: Backfitting; Change points; Dioxin; Kernel estimation; Median regression; Nonparametric regression; Partial linear model; Profile kernel methods; Semiparametric estimation; Undersmoothing.

Short Title: Backfitting and Profiling

¹Financial support from the IAP research network nr. P5/24 of the Belgian government (Belgian Science Policy) is gratefully acknowledged.

²Research supported by a grant from the National Cancer Institute (CA-57030), and by the Texas A&M Center for Environmental and Rural Health via a grant from the National Institute of Environmental Health Sciences (P30-ES09106).

1 Introduction

Consider a semiparametric problem that depends on a parameter β_0 and an unknown function $\theta_0(\cdot)$. The purpose of this paper is to compare backfitting and profiling methods in semiparametric regression. Our context is quite general and allows for estimation based on non-smooth criterion functions.

We first introduce the context of the general problem. Assume that the data (X_i, Y_i) ($i = 1, \dots, n$) are independent replications of a $(1 + d_y)$ -dimensional random vector (X, Y) . Let β denote a $q \times 1$ vector of parameters of interest, with true value β_0 , belonging to a compact subset \mathcal{B} of \mathbb{R}^q . Let $\theta = \theta(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ be an infinite dimensional ‘nuisance’ parameter with true value $\theta_0(\cdot)$, and let $\mathcal{L}\{Y, \beta_0, \theta_0(X)\}$ be a real-valued maximizing function for β_0 and θ_0 , in the sense that $E[\mathcal{L}_\beta\{Y, \beta_0, \theta_0(X)\}] = 0$ and $E[\mathcal{L}_\theta\{Y, \beta_0, \theta_0(x)\} | X = x] = 0$ for all x , where $\mathcal{L}_\beta\{y, \beta, \theta(x)\}$ denotes the vector of partial derivatives of $\mathcal{L}\{y, \beta, \theta(x)\}$ with respect to the components of β , and $\mathcal{L}_\theta\{y, \beta, \theta(x)\}$ denotes the partial derivative of $\mathcal{L}(y, \beta, z)$ with respect to z , and evaluated at $z = \theta(x)$. Inference for β_0 is then carried out by maximizing

$$n^{-1} \sum_{i=1}^n \mathcal{L}\{Y_i, \beta, \theta(X_i)\} \quad (1)$$

with respect to β for some $\theta(\cdot)$.

In the semiparametric literature, two approaches have been considered to maximize expression (1), these approaches differing in the way they treat the unknown function $\theta_0(\cdot)$.

The backfitting procedure has been investigated by many authors in special contexts, including Rice (1986), Speckman (1988), Buja et al. (1989), Hastie and Tibshirani (1990), Opsomer and Ruppert (1997, 1999), Mammen et al. (1999), Wand (1999) and Opsomer (2000). The basic idea is one of iteration. For any given β , let $\hat{\theta}(\cdot, \beta)$ be an estimate of $\theta_0(\cdot)$: the estimator we use is defined in the next section. Then define

$$m_{BF}\{y, \beta, \theta(x)\} = \mathcal{L}_\beta\{y, \beta, \theta(x)\}.$$

The backfitting estimator $\hat{\beta}_{BF}$ is now defined by the value of β that minimizes

$$\|n^{-1} \sum_{i=1}^n m_{BF}\{Y_i, \beta, \hat{\theta}(X_i, \beta)\}\| \quad (2)$$

over \mathcal{B} , where $\|\cdot\|$ is the Euclidean norm in \mathbb{R}^q .

The profile method also has a large literature, see for example Severini and Wong (1992), Severini and Staniswalis (1994), Carroll et al. (1997) and Murphy and van der Vaart (2000), among many others. This method again starts with $\hat{\theta}(x, \beta)$, which may be different from that constructed for backfitting, but it obtains the estimate of β differently. Specifically, it creates an estimating function m_{PR} by differentiating $\mathcal{L}\{y, \beta, \theta(x, \beta)\}$ with respect to β , i.e.,

$$\begin{aligned} m_{PR}\{y, \beta, \theta(x, \beta), \theta_\beta(x, \beta)\} &= \frac{d}{d\beta} \mathcal{L}\{y, \beta, \theta(x, \beta)\} \\ &= \mathcal{L}_\beta\{y, \beta, \theta(x, \beta)\} + \mathcal{L}_\theta\{y, \beta, \theta(x, \beta)\} \frac{\partial}{\partial \beta} \theta(x, \beta), \end{aligned}$$

where $\theta_\beta(x, \beta) = \frac{\partial}{\partial \beta} \theta(x, \beta)$ for any $\theta(x, \beta)$. With these definitions, the profile estimator $\hat{\beta}_{PR}$ is the value of β in \mathcal{B} for which

$$\left\| n^{-1} \sum_{i=1}^n m_{PR}\{Y_i, \beta, \hat{\theta}(X_i, \beta), \hat{\theta}_\beta(X_i, \beta)\} \right\| \quad (3)$$

is minimal, where $\hat{\theta}(\cdot, \beta)$ and $\hat{\theta}_\beta(\cdot, \beta)$ are defined in Section 2.

The comparison of backfitting and profiling has been the subject of some limited research. Consider a Gaussian model with independent data, scalar response \mathcal{Y}_i and predictors (Z_i, X_i) , so that in our context $Y_i = (\mathcal{Y}_i, Z_i)$, and suppose that the true mean is $Z_i^T \beta_0 + \theta_0(X_i)$. Opsomer and Ruppert (1999) showed that under certain conditions, backfitting and profiling produce asymptotically equivalent estimators, but only when backfitting an estimated function $\hat{\theta}(x, \beta)$ undersmoothed compared to that used by profiling. In more global contexts, with correlated data and multiple arguments for the function, backfitting and profiling are no longer necessarily asymptotically equivalent, see Hu, et al. (2004) for a counterexample.

In this note we study the two methods when the criterion functions $m_{BF}(y, \beta, z)$ and $m_{PR}(y, \beta, z_1, z_2)$ are not necessarily smooth in β and/or z , and when θ is estimated by kernel smoothing. We will prove that, under certain regularity conditions, the two methods are asymptotically equivalent, but only when the driving estimation methods $\hat{\theta}(x, \beta)$ employ different amounts of smoothness.

The paper is organized as follows. In the next section we introduce notation and develop the general conditions under which the estimators $\hat{\beta}_{BF}$ and $\hat{\beta}_{PR}$ are asymptotically normal.

We show that under certain primitive conditions, the profile estimator and the backfitting estimator have the same asymptotic variance. Section 3 gives two applications. The first, discussed in Section 3.1 deals with an application of the general theory to partial linear median regression. The second, in Section 3.2, is concerned with a varying coefficient change point model, motivated by a problem in toxicology. The proofs and the conditions under which the main results for backfitting and profiling are valid are given in Appendix A and B, respectively.

2 Main Results

Let $K_h(u) = h^{-1}K(u/h)$, K be a symmetric kernel density function and h be a smoothing parameter. For the backfitting procedure, let $\hat{\theta}(x, \beta)$ be defined by a value of θ that maximizes

$$n^{-1} \sum_{i=1}^n K_h(X_i - x) \mathcal{L}(Y_i, \beta, \theta), \quad (4)$$

for fixed values of β and x . In order to focus on the primary issues, we assume the existence of a well-defined maximizer of (4).

For the profiling estimator, all we need is that $\hat{\theta}(x, \beta)$ and $\hat{\theta}_\beta(x, \beta)$ satisfy assumption (PR1) given in Appendix B. This implies that the asymptotic distribution of $\hat{\beta}$ does not depend on the estimators of θ_0 and $\theta_{0\beta}$, as long as assumption (PR1) is fulfilled. While other nonparametric estimators based on e.g. splines or local polynomials can be used, in particular, $\hat{\theta}$ and $\hat{\theta}_\beta$ can be estimated in the following way: let $\hat{\theta}(x, \beta)$ be defined as for the backfitting procedure and let $\hat{\theta}_\beta(x, \beta)$ be the partial derivative of $\hat{\theta}(x, \beta)$ with respect to β , or, in case $\hat{\theta}(x, \beta)$ is not differentiable with respect to β , define $\hat{\theta}_\beta(x, \beta)$ by

$$\frac{\partial}{\partial \beta} \int \hat{\theta}(x, b) L_g(\beta - b) db$$

where L is a kernel density function and g is an appropriate bandwidth. Recall that $\hat{\beta}_{BF}$ and $\hat{\beta}_{PR}$ are the estimators of β_0 defined in (2) and (3).

Let $\theta_0(x, \beta)$ denote a solution of $E\{\mathcal{L}_\theta(Y, \beta, \theta) | X = x\} = 0$ with respect to θ for fixed β and x , where the expectation is calculated under the distribution induced by $\{\beta_0, \theta_0(\cdot)\}$. We assume that $\theta_0(x, \beta)$ is unique. Clearly, $\theta_0(\cdot, \beta_0) \equiv \theta_0(\cdot)$.

For any function $H \equiv \{H_1(\beta, \theta), \dots, H_d(\beta, \theta)\}$ of (say) dimension d , we use the notation $\frac{\partial}{\partial \beta} H(\beta, \theta)$ and $\frac{d}{d\beta} H(\beta, \theta)$ to denote the $d \times q$ matrix with (i, j) th-element ($i = 1, \dots, d, j = 1, \dots, q$) given by

$$\frac{\partial}{\partial \beta} H(\beta, \theta)_{ij} = \lim_{\tau \rightarrow 0} \frac{1}{\tau} [H_i\{\beta + \tau e_j, \theta(\cdot, \beta)\} - H_i\{\beta, \theta(\cdot, \beta)\}] \quad (5)$$

and

$$\frac{d}{d\beta} H(\beta, \theta)_{ij} = \lim_{\tau \rightarrow 0} \frac{1}{\tau} [H_i\{\beta + \tau e_j, \theta(\cdot, \beta + \tau e_j)\} - H_i\{\beta, \theta(\cdot, \beta)\}] \quad (6)$$

respectively, where $e_j = (e_{j1}, \dots, e_{jq})$ and $e_{jk} = \delta_{jk} = I(j = k)$ ($k = 1, \dots, q$).

We are now ready to state the main result concerning the asymptotic normality of $\widehat{\beta}_{BF}$ and $\widehat{\beta}_{PR}$. Let $\mathcal{G}(\beta) = \frac{d}{d\beta} E[\mathcal{L}_\beta\{Y, \beta, \theta_0(X, \beta)\}]$ and define

$$\Sigma = \text{cov}\left[\mathcal{L}_\beta\{Y, \beta_0, \theta_0(X, \beta_0)\} + \mathcal{L}_\theta\{Y, \beta_0, \theta_0(X, \beta_0)\} \frac{\partial}{\partial \beta} \theta_0(X, \beta_0)\right].$$

Theorem 2.1 *Assume (BF1)–(BF8): in particular, assume that the bandwidth h satisfies $nh^4 \rightarrow 0$, and not the usual optimal bandwidth of $h \propto n^{-1/5}$. Then,*

$$n^{1/2}(\widehat{\beta}_{BF} - \beta_0) \xrightarrow{d} \text{Normal}\{0, \mathcal{G}^{-1}(\beta_0)\Sigma\mathcal{G}^{-1}(\beta_0)^T\}.$$

Theorem 2.2 *Assume (PR1)–(PR4): in particular, we allow that the bandwidth h satisfies $h \propto n^{-1/5}$. Then,*

$$n^{1/2}(\widehat{\beta}_{PR} - \beta_0) \xrightarrow{d} \text{Normal}\{0, \mathcal{G}^{-1}(\beta_0)\Sigma\mathcal{G}^{-1}(\beta_0)^T\}.$$

The proof of these results, as well the assumptions under which they are valid, can be found in Appendix A for the backfitting method and in Appendix B for the profiling method.

As a consequence, the backfitting and profiling method produce asymptotically equivalent estimators, also when m_{BF} or m_{PR} are not smooth in β or θ . The backfitting procedure requires however that undersmoothing be used to estimate $\widehat{\theta}(x, \beta)$, whereas the profiling procedure does not.

Note that, although the asymptotic variance $\mathcal{G}^{-1}(\beta_0)\Sigma\mathcal{G}^{-1}(\beta_0)^T$ has an explicit formula, its actual computation might be complicated in certain situations. In such cases, a bootstrap approximation can be useful. See Theorem B in Chen, Linton and Van Keilegom (2003) for general conditions under which a naive bootstrap procedure is valid.

3 Applications

3.1 Partially Linear Median Regression

Consider the model

$$\mathcal{Y}_i = Z_i^T \beta_0 + \theta_0(X_i) + \epsilon_i \quad (7)$$

where Z_i is a possibly vector-valued covariate of dimension q , X_i is a scalar covariate, and $\text{med}(\epsilon_i | X_i, Z_i) = 0$. In our general notation, let $Y = (\mathcal{Y}, Z)$. The criterion function is

$$\mathcal{L}\{Y, \beta, \theta(X)\} = -|\mathcal{Y} - Z^T \beta - \theta(X)|. \quad (8)$$

The backfitting estimator $\hat{\beta}_{BF}$ of β_0 has been considered in Chen, Linton and Van Keilegom (2003), see their Example 2. Here, we consider the profiling estimator.

It is readily seen that for fixed β , $\theta_0(x, \beta) = \text{med}(\mathcal{Y} - Z^T \beta | X = x)$. Let $\hat{\theta}(x, \beta)$ be the kernel estimator of the conditional median of $\mathcal{Y} - Z^T \beta$ given $X = x$. Note that $\hat{\theta}(x, \beta)$ is not smooth in β , because $\hat{\theta}(x, \beta)$ is piecewise constant as a function of β . Hence we define $\hat{\theta}_\beta(x, \beta)$ by

$$\hat{\theta}_\beta(x, \beta) = \frac{\partial}{\partial \beta} \int \hat{\theta}(x, b) L_g(\beta - b) db.$$

Let $\Theta = \{\theta : \theta(\cdot, \beta) \in C_M^\alpha(R_X) \text{ for all } \beta\}$ for some $\alpha > 1$, $0 < M < \infty$, and some compact interval R_X , see van der Vaart and Wellner (1996), p. 154 for the definition of the class $C_M^\alpha(R_X)$. Assume that θ_0 and the components of $\theta_{0\beta}$ belong to Θ . Then, using kernel theory for median regression (see e.g. Chaudhuri (1991)), it can be seen that assumption (PR1) is valid : for $\hat{\theta}_\beta - \theta_{0\beta}$ note that

$$\begin{aligned} \|\hat{\theta}_\beta - \theta_{0\beta}\|_\infty &\leq \left\| \frac{\partial}{\partial \beta} \int \{\hat{\theta}(x, b) - \theta_0(x, b)\} L_g(\beta - b) db \right\|_\infty \\ &\quad + \left\| \frac{\partial}{\partial \beta} \int \{\theta_0(x, b) - \theta_0(x, \beta)\} L_g(\beta - b) db \right\|_\infty \end{aligned}$$

and this is $o_P(n^{-1/4})$ provided L is a symmetric, compactly supported kernel function, $ng^8 \rightarrow 0$, $nh^2g^4 \rightarrow \infty$, $nh^8g^{-4} \rightarrow 0$, $\|\hat{\theta} - \theta_0\|_\infty = O_P\{(nh)^{-1/2} + h^2\}$ and θ_0 is three times continuously differentiable with respect to the components of β . For example, take $h = C_1 n^{-1/5}$ and $g = C_2 n^{-1/7}$ for some $C_1, C_2 > 0$. In order to show that $\|\hat{\theta} - \theta_0\|_\infty = O_P\{(nh)^{-1/2} + h^2\}$,

note that Chaudhuri (1991) shows that $\sup_x |\hat{\theta}(x, \beta) - \theta_0(x)| = O_P\{(nh)^{-1/2} + h^2\}$ for fixed β and for $h \propto n^{-1/5}$. It is possible to extend this result to bandwidths h that satisfy the above constraints and to prove that the given rate holds uniformly over all β . It suffices for this to replace the supremum over β in an appropriate way by a maximum over a set of grid points (of size tending to infinity) and to prove the consistency uniformly over the set of grid points.

Direct calculations show that

$$\begin{aligned} E[\mathcal{L}_\beta\{Y, \beta, \theta(X, \beta)\}|X] &= -E([2F_{Y|X,Z}\{Z^T\beta + \theta(X, \beta)\} - 1]Z|X), \\ E[\mathcal{L}_\theta\{Y, \beta, \theta(X, \beta)\}|X] &= -E[2F_{Y|X,Z}\{Z^T\beta + \theta(X, \beta)\} - 1|X]. \end{aligned}$$

In addition,

$$\mathcal{G}(\beta_0) = -2E\left[f_{Y|X,Z}\{Z^T\beta_0 + \theta_0(X)\}Z\left\{Z + \frac{\partial}{\partial\beta}\theta_0(X, \beta_0)\right\}^T\right].$$

Hence, assumption (PR2) is verified under standard smoothness conditions on $F_{Y|X,Z}$. Also, (PR4) holds under classical identifiability conditions. It is readily seen that (PR3)(i) is valid for $r_\ell = 2$ and $s_\ell = 1/2$. Finally, for assumption (PR3)(ii) we make use of Theorem 2.7.1 in van der Vaart and Wellner (1996). It is easily checked that

$$\int_0^\infty \sqrt{\log N(\varepsilon^2, \tilde{\Theta}, \|\cdot\|_\infty)} d\varepsilon \leq C \int_0^{(2M)^{1/2}} \varepsilon^{-1/\alpha} d\varepsilon < \infty,$$

for some $C > 0$. The asymptotic normality of $\hat{\beta}_{PR}$ now follows. Note that the matrix Σ equals

$$\Sigma = \text{cov}\left(\{2I(\epsilon \geq 0) - 1\}\left[Z - \frac{E\{f_{\epsilon|X,Z}(0)Z|X\}}{f_{\epsilon|X}(0)}\right]\right),$$

since it is easily seen that $\theta_{0\beta}(X, \beta_0) = -E\{f_{\epsilon|X,Z}(0)Z|X\}/f_{\epsilon|X}(0)$.

3.2 Varying Coefficient Change Point Model

Consider the following model

$$\mathcal{Y}_i = \theta_{01}(X_i) + \theta_{02}(X_i)|Z_i - \beta_0|_+ + \epsilon_i, \tag{9}$$

where $E(\epsilon_i|X_i, Z_i) = 0$, X_i and Z_i are scalar covariates and $z_+ = zI(z > 0)$. An interesting application of this type of model can be found in toxicology, where models of the form

$$E(\mathcal{Y}_i|Z_i) = \theta_{01} + \theta_{02}|Z_i - \beta_0|_+^{\lambda_0} \quad (10)$$

are compatible with accepted understanding of the basic structure of dose-response curves for exposure to dioxin. Roberts (1991) states that “new findings suggest that responses to dioxin increase slowly at first but then shoot up after passing a critical concentration”. Indeed, researchers have “agreed that before dioxin can cause any of its myriad toxic effects . . . it must first bind to and then activate a receptor. . . . If receptor binding is indeed the essential first step . . . then that implies there is a safe dose or practical threshold below which no toxic effects occur”. Feder (1975) constructs \sqrt{n} -consistent and asymptotically normally distributed estimators for $(\theta_{01}, \theta_{02}, \beta_0)$ in model (10) when $\lambda_0 = 1$, see his Example 1 on p. 77, setting his $\theta_{12} = 0$. Model (9) goes one step further, in the sense that it allows the average response before and the slope after critical concentration to depend on e.g. age or any other individual characteristic.

Let $Y = (\mathcal{Y}, Z)$ and define the following criterion function :

$$\mathcal{L}\{Y, \beta, \theta(X)\} = -\{\mathcal{Y} - \theta_1(X) - \theta_2(X)|Z - \beta|_+\}^2.$$

Note that the theory developed in Section 2 can be extended in an obvious way to bivariate nuisance functions. Straightforward calculations show that for fixed β ,

$$\begin{aligned} \theta_{02}(X, \beta) &= \theta_{02}(X) \frac{\text{Cov}(|Z - \beta_0|_+, |Z - \beta|_+|X)}{\text{Var}(|Z - \beta|_+|X)} = \frac{\text{Cov}(\mathcal{Y}, |Z - \beta|_+|X)}{\text{Var}(|Z - \beta|_+|X)}, \\ \theta_{01}(X, \beta) &= E(\mathcal{Y}|X) - \theta_{02}(X, \beta)E(|Z - \beta|_+|X). \end{aligned}$$

Also, let $\hat{\theta}_1(X, \beta)$ and $\hat{\theta}_2(X, \beta)$ be the estimators obtained by replacing the conditional means, variances and covariances in the above expressions by the corresponding kernel estimators, and let $\hat{\theta}_{1\beta}(X, \beta)$ and $\hat{\theta}_{2\beta}(X, \beta)$ be obtained by replacing $|Z_i - \beta|_+$ in these kernel estimators by $-I(Z_i \geq \beta)$ ($i = 1, \dots, n$).

As for the example on partial linear median regression, the main assumptions to verify are (BF7) for the backfitting procedure and (PR1) and (PR3) for the profiling method. We start with (PR1). Let $\Theta = \{\theta : \theta(\cdot, \beta) \in C_M^\alpha(R_X) \text{ for all } \beta\}$ for some $\alpha > 1/2$, $0 < M < \infty$

and some compact interval R_X , see van der Vaart and Wellner (1996), page 154. Assume that θ_0 and $\theta_{0\beta}$ belong to Θ . We will show that $\|\widehat{\theta}_{j\beta} - \theta_{0j\beta}\|_\infty = o_P(n^{-1/4})$ ($j = 1, 2$): the other conditions in (PR1) can be proved similarly. Since $\theta_{0j}(\cdot, \beta)$ is composed of variances, covariances and means, it suffices to consider each of these factors separately. For simplicity, we restrict attention to the mean, i.e. we consider

$$\sup_{x, \beta} \left| n^{-1} \sum_{i=1}^n \frac{K_h(X_i - x)}{\sum_{j=1}^n K_h(X_j - x)} I(Z_i \geq \beta) - P(Z \geq \beta | X = x) \right| = o_P(n^{-1/4})$$

provided $nh^2 \rightarrow \infty$ and $nh^8 \rightarrow 0$, see e.g. Proposition 4.1 in Akritas and Van Keilegom (2001).

Part (i) of conditions (BF7) and (PR3) can be easily seen to hold true for $s_\ell = 1$ and $r_\ell = 2$. For part (ii), since $N(\varepsilon, C_M^\alpha(R_X), \|\cdot\|_\infty) = O\{\exp(K\varepsilon^{-1/\alpha})\}$, see Theorem 2.7.1 in van der Vaart and Wellner (1996), it follows that the integral in part (ii) is finite. The asymptotic normality of both $\widehat{\beta}_{BF}$ and $\widehat{\beta}_{PR}$ now follows. The calculation of the asymptotic variance is straightforward but leads to lengthy formulas, and is left to the reader.

Appendix A: Proofs for Backfitting

Make the definitions

$$\begin{aligned} M_{nBF}(\beta, \theta) &= n^{-1} \sum_{i=1}^n m_{BF}\{Y_i, \beta, \theta(X_i, \beta)\}, \\ M_{BF}(\beta, \theta) &= E[m_{BF}\{Y, \beta, \theta(X, \beta)\}], \end{aligned}$$

and define the $q \times q$ matrix $\Gamma_{BF, \beta}(\beta, \theta) = \frac{d}{d\beta} M_{BF}(\beta, \theta) = \frac{d}{d\beta} E[\mathcal{L}_\beta\{Y, \beta, \theta(X, \beta)\}]$. Also, for a function $\xi(\cdot) = \xi(X, \beta)$, let $\Gamma_{BF, \theta}(\beta, \theta)[\xi]$ denote the Gâteaux-derivative of $M_{BF}(\beta, \theta)$ in the direction ξ , i.e.,

$$\begin{aligned} \Gamma_{BF, \theta}(\beta, \theta)[\xi] &= \lim_{\tau \rightarrow 0} \frac{1}{\tau} \{M_{BF}(\beta, \theta + \tau\xi) - M_{BF}(\beta, \theta)\} \\ &= \lim_{\tau \rightarrow 0} \frac{1}{\tau} E[\mathcal{L}_\beta\{Y, \beta, (\theta + \tau\xi)(X, \beta)\} - \mathcal{L}_\beta\{Y, \beta, \theta(X, \beta)\}] \\ &= E\left(\frac{\partial}{\partial \theta} E[\mathcal{L}_\beta\{Y, \beta, \theta(X, \beta)\} | X] \xi(X, \beta)\right), \end{aligned}$$

where $\frac{\partial}{\partial \theta} E[\mathcal{L}_\beta\{Y, \beta, \theta(X, \beta)\} | X] = \frac{\partial}{\partial z} E[\mathcal{L}_\beta(Y, \beta, z) | X]_{z=\theta(X, \beta)}$. Note that $M_{BF}(\beta_0, \theta_0) = 0$.

For any function $g = (g_1, \dots, g_d)$ of (say) dimension d defined on a set \mathcal{A} in \mathbb{R}^a , for any $y \in \mathcal{A}$ and any k , let $\frac{\partial^k}{\partial y^k} g(y)$ denote the vector of all partial derivatives of order k of the form $\frac{\partial^k}{\partial y_1^{k_1} \dots \partial y_a^{k_a}} g_j(y)$, where $\sum_{i=1}^a k_i = k$ and $1 \leq j \leq d$. Let $\|g\|_\infty = \max_{1 \leq j \leq d} \sup_{y \in \mathcal{A}} |g_j(y)|$. In particular, for a function $\theta = \theta(x, \beta)$, $\|\theta\|_\infty = \sup_{x, \beta} |\theta(x, \beta)|$, and $\|\frac{\partial \theta}{\partial \beta}\|_\infty = \max_{1 \leq \ell \leq q} \sup_{x, \beta} |\frac{\partial \theta}{\partial \beta_\ell}(x, \beta)|$.

Further, let Θ be some space of functions $\theta = \theta(x, \beta)$ ($x \in \mathbb{R}$, $\beta \in \mathcal{B}$) for which $\|\theta\|_\infty \leq M$ for some $M > 0$.

The conditions below use the concept of covering number which is defined as follows. For $\epsilon > 0$ and any normed space $(\Theta, \|\cdot\|)$ of functions, the covering number $N(\epsilon, \Theta, \|\cdot\|)$ is the minimal number of balls $\{\eta : \|\eta - \theta\| < \epsilon\}$ of radius ϵ needed to cover Θ . The centers of the balls need not belong to Θ , but they should have finite norms.

(BF1) The bandwidth h satisfies $nh^4 \rightarrow 0$ as n tends to infinity.

(BF2) The probability density function K has compact support and $\int uK(u) du = 0$.

(BF3) X is absolutely continuous and has compact support R_X , its density f_X is twice continuously differentiable and $\inf_x f_X(x) > 0$.

(BF4) $\theta_0 \in \Theta$, $\frac{\partial^{k+l}}{\partial x^k \partial \beta^\ell} \theta_0(x, \beta)$ ($0 \leq k+l \leq 3$) exists for almost all x and β and $\|\frac{\partial^{k+l} \theta_0}{\partial x^k \partial \beta^\ell}\|_\infty < \infty$.

(BF5) (i) $P(\hat{\theta} \in \Theta) \rightarrow 1$ as $n \rightarrow \infty$ and $\|\hat{\theta} - \theta_0\|_\infty = o_P(n^{-1/4})$.

(ii) $\sup_x |(\hat{\theta} - \theta_0)(x, \hat{\beta}) - (\hat{\theta} - \theta_0)(x, \beta_0)| = o_P(1) \|\hat{\beta} - \beta_0\|$.

(iii) $\sup_x |n^{-1} \sum_{i=1}^n K_h(X_i - x) \mathcal{L}_\theta\{Y_i, \beta_0, \hat{\theta}(x, \beta_0)\}| = o_P(n^{-1/2})$.

(BF6) (i) For all y , $\mathcal{L}(y, \beta, \theta)$ is differentiable with respect to β and θ , for almost all β and θ .

(ii) $\frac{\partial}{\partial \theta} E[\mathcal{L}_\beta\{Y, \beta, \theta_0(X, \beta)\}|X]$ and $\frac{\partial}{\partial \beta} E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}|X]$ exist for all $\beta \in \mathcal{B}$, and they are equal.

(iii) $E\left\{\sup_{|\theta| \leq M} |\mathcal{L}_\theta(Y, \beta_0, \theta)|^2\right\} < \infty$.

(iv) $\frac{\partial^{j+k+\ell}}{\partial \theta^j \partial x^k \partial \beta^\ell} E\{\mathcal{L}_\beta(Y, \beta, \theta)|X = x\}$ and $\frac{\partial^{j+k+\ell}}{\partial \theta^j \partial x^k \partial \beta^\ell} E\{\mathcal{L}_\theta(Y, \beta, \theta)|X = x\}$ exist for $0 \leq j+k+\ell \leq 2$ and for all β, θ and x , and

$$\sup_{\beta \in \mathcal{B}, |\theta| \leq M, x \in R_X} \left| \frac{\partial^{j+k+\ell}}{\partial \theta^j \partial x^k \partial \beta^\ell} E\{\mathcal{L}_\beta(Y, \beta, \theta)|X = x\} \right| < \infty,$$

$$\sup_{\beta \in \mathcal{B}, |\theta| \leq M, x \in R_X} \left| \frac{\partial^{j+k+\ell}}{\partial \theta^j \partial x^k \partial \beta^\ell} E\{\mathcal{L}_\theta(Y, \beta, \theta) | X = x\} \right| < \infty,$$

(v) $\mathcal{G}(\beta)$ exists for β in a neighborhood of β_0 , is continuous at β_0 and $\mathcal{G}(\beta_0)$ is of full rank.

(BF7) (i)

$$E \left\{ \sup_{(\beta', \theta') : \|\beta - \beta'\| \leq \delta, \|\theta - \theta'\|_\infty \leq \delta} |\mathcal{L}_\theta(Y, \beta, \theta) - \mathcal{L}_\theta(Y, \beta', \theta')|^{r_0} \right\} \leq K_0 \delta^{r_0 s_0},$$

$$E \left\{ \sup_{(\beta', \theta') : \|\beta - \beta'\| \leq \delta, \|\theta - \theta'\|_\infty \leq \delta} |\mathcal{L}_{\beta, \ell}(Y, \beta, \theta) - \mathcal{L}_{\beta, \ell}(Y, \beta', \theta')|^{r_\ell} \right\} \leq K_\ell \delta^{r_\ell s_\ell},$$

for $r_0 = 2, 2 + \eta$, for some $r_\ell \geq 2$ ($\ell = 1, \dots, q$), for all $(\beta, \theta) \in \mathcal{B} \times \Theta$, all $\delta > 0$, for some $\eta > 0$, some $0 < s_\ell \leq 1$ and some $K_\ell > 0$ ($\ell = 0, \dots, q$).

(ii) $\int_0^\infty \sqrt{\log N(\varepsilon^{1/s_\ell}, \tilde{\Theta}, \|\cdot\|_\infty)} d\varepsilon < \infty$, for $\ell = 0, \dots, q$, where $\tilde{\Theta} = \{\theta(\cdot, \beta) : \theta \in \Theta, \beta \in \mathcal{B}\}$.

(BF8)(i) For all $\delta > 0$, there exists a $\varepsilon > 0$ such that $\inf_{\|\beta - \beta_0\| > \delta} \|M_{BF}(\beta, \theta_0)\| \geq \varepsilon$.

(ii) Uniformly for all $\beta \in \mathcal{B}$, $M_{BF}(\beta, \theta)$ is continuous in θ at θ_0 (with respect to the $\|\cdot\|_\infty$ norm).

(iii) $\Gamma_{BF, \theta}(\beta, \theta_0)[\theta - \theta_0]$ exists in all directions $\theta - \theta_0 \in \Theta$.

Assumption (BF1) requires that undersmoothing is used to estimate the nuisance function θ_0 . This will not be required for the profiling method. Conditions (BF2)-(BF6) and (BF8) are standard regularity conditions that can be easily verified in practical situations. Note that assumption (BF6) does not impose smoothness conditions on \mathcal{L}_β and \mathcal{L}_θ , but requires instead that $E(\mathcal{L}_\beta)$ and $E(\mathcal{L}_\theta)$ are differentiable.

The condition on the covering number in (BF7) can be checked by using e.g. the results obtained by van der Vaart and Wellner (1996). A common special case is the case where the class $\tilde{\Theta}$ belongs to $C_M^\alpha(R_X)$, which is a certain subclass of the space of all functions that possess partial derivatives up to order $\alpha > 0$, see page 154 in their book for a precise definition. Theorem 2.7.1 (page 155) gives a bound on the covering number for this space.

For the proofs below, we restrict attention for simplicity to the case $q = 1$. The general case $q \geq 1$ can be obtained in a similar way, but requires more complex notation.

We start with a technical lemma.

Lemma A.1 Assume (BF1)–(BF8). Then,

$$\begin{aligned} & n^{-1} \sum_{i=1}^n E_X \left(\frac{K_h(X_i - X)}{f_X(X)} \frac{\partial}{\partial \beta} \theta_0(X, \beta_0) [\mathcal{L}_\theta\{Y_i, \beta_0, \theta_0(X_i)\} - \mathcal{L}_\theta\{Y_i, \beta_0, \hat{\theta}(X, \beta_0)\}] \right) \\ &= E_{X_1, X_2, Y_1} \left(\frac{K_h(X_1 - X_2)}{f_X(X_2)} \frac{\partial}{\partial \beta} \theta_0(X_2, \beta_0) [\mathcal{L}_\theta\{Y_1, \beta_0, \theta_0(X_1)\} - \mathcal{L}_\theta\{Y_1, \beta_0, \hat{\theta}(X_2, \beta_0)\}] \right) \\ &\quad + o_P(n^{-1/2}), \end{aligned}$$

where the expectations are taken conditionally on the data (X_i, Y_i) ($i = 1, \dots, n$).

Proof. Throughout the proof, C denotes a generic constant, whose value may change from one line to another. The following abbreviated notations will be used : let $\mathcal{H}(Y, \theta) = \mathcal{L}_\theta(Y, \beta_0, \theta)$ and $g(X) = \frac{\partial}{\partial \beta} \theta_0(X, \beta_0)$. To prove this result we will make use of modern empirical process theory see e.g. van der Vaart and Wellner 1996. Consider the process $\sum_{i=1}^n Z_{ni}(\theta)$, where

$$\begin{aligned} Z_{ni}(\theta) &= n^{-1/2} \left\{ E_X \left(\frac{K_h(X_i - X)}{f_X(X)} g(X) [\mathcal{H}\{Y_i, \theta_0(X_i)\} - \mathcal{H}\{Y_i, \theta(X)\}] \right) \right. \\ &\quad \left. - E_{X_1, X_2, Y_1} \left(\frac{K_h(X_1 - X_2)}{f_X(X_2)} g(X_2) [\mathcal{H}\{Y_1, \theta_0(X_1)\} - \mathcal{H}\{Y_1, \theta(X_2)\}] \right) \right\}, \end{aligned}$$

where θ belongs to Θ . For simplicity we suppress the dependence of θ on β_0 . Note that by assumption (BF5)(i), $P(\hat{\theta} \in \Theta) \rightarrow 1$. In order to show the weak convergence of this process we will verify the conditions of Theorem 2.11.9 in van der Vaart and Wellner (1996):

$$\sum_{i=1}^n E \left[\sup_{\theta \in \Theta} |Z_{ni}(\theta)| I \left\{ \sup_{\theta \in \Theta} |Z_{ni}(\theta)| > \eta \right\} \right] \rightarrow 0 \quad \text{for every } \eta > 0; \quad (11)$$

$$\int_0^{\delta_n} \sqrt{\log N_{[]}(\varepsilon, \Theta, L_2^n)} d\varepsilon \rightarrow 0 \quad \text{for every } \delta_n \downarrow 0; \quad (12)$$

$$\sum_{i=1}^n Z_{ni}(\theta) \text{ converges marginally for every } \theta \in \Theta, \quad (13)$$

where $N_{[]}(\varepsilon, \Theta, L_2^n)$ is the bracketing number, defined as the minimal number of sets N_ε in a partition $\Theta = \cup_{j=1}^{N_\varepsilon} \Theta_{\varepsilon j}$, such that for every $j = 1, \dots, N_\varepsilon$:

$$\sum_{i=1}^n E \left\{ \sup_{\theta_1, \theta_2 \in \Theta_{\varepsilon j}} |Z_{ni}(\theta_1) - Z_{ni}(\theta_2)|^2 \right\} \leq \varepsilon^2. \quad (14)$$

The first two conditions (11) and (12) imply the asymptotic tightness of the process and can be proved separately for the four terms in the definition of $\sum_{i=1}^n Z_{ni}$. We will restrict

ourselves to showing (11) and (12) for the second term:

$$\sum_{i=1}^n \tilde{Z}_{ni}(\theta) = n^{-1/2} \sum_{i=1}^n E_X \left[\frac{K_h(X_i - X)}{f_X(X)} g(X) \mathcal{H}\{Y_i, \theta(X)\} \right].$$

We start with verifying condition (12). Fix $\varepsilon > 0$. From assumption (BF7)(ii) it follows that there exist functions $\theta_1, \dots, \theta_{N_\varepsilon}$ in Θ such that $\int_0^{\delta_n} \sqrt{\log N_\varepsilon} d\varepsilon \rightarrow 0$ and such that the balls $\{\theta : \|\theta - \theta_j\|_\infty \leq \varepsilon^{1/s_0}\}$ ($j = 1, \dots, N_\varepsilon$) cover Θ . We will show that for any $1 \leq j \leq N_\varepsilon$,

$$\sum_{i=1}^n E \left\{ \sup_{\|\theta - \theta_j\|_\infty \leq \varepsilon^{1/s_0}} |\tilde{Z}_{ni}(\theta) - \tilde{Z}_{ni}(\theta_j)|^2 \right\} \leq \varepsilon^2. \quad (15)$$

The left hand side of (15) equals

$$\begin{aligned} & E \left(\sup_{\|\theta - \theta_j\|_\infty \leq \varepsilon^{1/s_0}} \left| \int K_h(X_1 - x) g(x) [\mathcal{H}\{Y_1, \theta(x)\} - \mathcal{H}\{Y_1, \theta_j(x)\}] dx \right|^2 \right) \\ &= E \left(\sup_{\|\theta - \theta_j\|_\infty \leq \varepsilon^{1/s_0}} \left| \int K(u) g(X_1 - hu) [\mathcal{H}\{Y_1, \theta(X_1 - hu)\} - \mathcal{H}\{Y_1, \theta_j(X_1 - hu)\}] du \right|^2 \right) \\ &\leq \sup_x |g(x)|^2 \int K(u) E \left[\sup_{\|\theta - \theta_j\|_\infty \leq \varepsilon^{1/s_0}} |\mathcal{H}\{Y_1, \theta(X_1 - hu)\} - \mathcal{H}\{Y_1, \theta_j(X_1 - hu)\}|^2 \right] du \\ &\leq C \sup_x |g(x)|^2 \varepsilon^2, \end{aligned}$$

where the last inequality follows from assumption (BF7)(i). This shows (15), up to a universal constant, and hence, (12) is satisfied for the class $\sum_{i=1}^n \tilde{Z}_{ni}(\theta)$. We next verify (11). With Z_{ni} replaced by \tilde{Z}_{ni} , the left hand side of (11) is bounded by

$$\begin{aligned} & n^{1/2} \sup_x |g(x)| E \left(\sup_\theta |\mathcal{L}_\theta(Y, \beta_0, \theta)| I \left[\sup_\theta |\mathcal{L}_\theta(Y, \beta_0, \theta)| > \eta n^{1/2} \left\{ \sup_x |g(x)| \right\}^{-1} \right] \right) \\ &= o(1), \end{aligned}$$

where we have used assumption (BF6)(iii). For the convergence of the marginals of $\sum_{i=1}^n Z_{ni}(\theta)$, we verify Liapunov's condition :

$$\frac{\sum_{i=1}^n E |Z_{ni}(\theta)|^{2+\eta}}{\left[\sum_{i=1}^n \text{Var}\{Z_{ni}(\theta)\} \right]^{(2+\eta)/2}} \rightarrow 0$$

for some $\eta > 0$. First, consider the variance. Using a similar derivation as above, we obtain for any $\theta \in \Theta$:

$$\sum_{i=1}^n \text{Var}\{Z_{ni}(\theta)\} \leq \sup_x |g(x)|^2 \int K(u) E |\mathcal{H}\{Y_1, \theta_0(X_1)\} - \mathcal{H}\{Y_1, \theta(X_1 - hu)\}|^2 du$$

$$\begin{aligned}
&\leq C \sup_x |g(x)|^2 \int K(u) \sup_x |\theta_0(x) - \theta(x - hu)|^{2s_0} du \\
&\leq Ch^{2s_0} \sup_x |g(x)|^2 \sup_x \left| \frac{\partial}{\partial x} \theta_0(x) \right|^{2s_0} \int K(u) |u|^{2s_0} du + 2CM \sup_x |g(x)|^2 \\
&= O(1).
\end{aligned} \tag{16}$$

In a similar way one can show that $\sum_{i=1}^n E|Z_{ni}(\theta)|^{2+\eta} = O(n^{-\eta/2})$, since assumption (BF7)(i) assures that $E|\mathcal{H}(Y_1, \theta_0(X_1)) - \mathcal{H}(Y_1, \theta(X_1 - hu))|^{2+\eta} \leq C \sup_x |\theta_0(x) - \theta(x - hu)|^{2s_0}$. Hence, the Liapunov ratio is $O\{n^{-\eta/2}\} = o(1)$. This shows the weak convergence of the process $\sum_{i=1}^n Z_{ni}(\theta)$ ($\theta \in \Theta$). It now follows that $\sup_{\theta \in \Theta} |\sum_{i=1}^n Z_{ni}(\theta)| = O_P(1)$. Finally, arguments similar to those in (16) show that $\sum_{i=1}^n \text{Var}\{Z_{ni}(\hat{\theta})\} = o_P(1)$ (where the variance is calculated conditionally on the value of $\hat{\theta}$), so that $\sum_{i=1}^n Z_{ni}(\hat{\theta}) = o_P(1)$, from which the result follows.

Lemma A.2 *Assume (BF1)–(BF8). Then,*

$$\Gamma_{BF,\theta}(\beta_0, \theta_0)[\hat{\theta} - \theta_0] = n^{-1} \sum_{i=1}^n \mathcal{L}_\theta\{Y_i, \beta_0, \theta_0(X_i, \beta_0)\} \frac{\partial}{\partial \beta} \theta_0(X_i, \beta_0) + o_P(n^{-1/2}). \tag{17}$$

Proof. Recall the definitions of $\frac{\partial}{\partial \beta}$ and $\frac{d}{d\beta}$ given in (5) and (6). First note that

$$\begin{aligned}
&\Gamma_{BF,\theta}(\beta_0, \theta_0)[\hat{\theta} - \theta_0] \\
&= \lim_{\tau \rightarrow 0} \frac{1}{\tau} E(\mathcal{L}_\beta[Y, \beta_0, \{\theta_0 + \tau(\hat{\theta} - \theta_0)\}](X, \beta_0) - \mathcal{L}_\beta\{Y, \beta_0, \theta_0(X, \beta_0)\}) \\
&= E\left(\frac{\partial}{\partial \theta} E[\mathcal{L}_\beta\{Y, \beta_0, \theta_0(X, \beta_0)\} | X](\hat{\theta} - \theta_0)(X, \beta_0)\right) \\
&= E\left(\frac{\partial}{\partial \beta} E[\mathcal{L}_\theta\{Y, \beta_0, \theta_0(X, \beta_0)\} | X](\hat{\theta} - \theta_0)(X, \beta_0)\right) \\
&= -E\left(\frac{\partial}{\partial \theta} E[\mathcal{L}_\theta\{Y, \beta_0, \theta_0(X, \beta_0)\} | X](\hat{\theta} - \theta_0)(X, \beta_0) \frac{\partial}{\partial \beta} \theta_0(X, \beta_0)\right),
\end{aligned} \tag{18}$$

since $E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\} | X] = 0$ for all β . Next, let $g(X) = \frac{\partial}{\partial \beta} \theta_0(X, \beta_0)$ and $\mathcal{H}(Y, \theta) = \mathcal{L}_\theta(Y, \beta_0, \theta)$. The right hand side of (17) equals

$$\begin{aligned}
&n^{-1} \sum_{i=1}^n E_X \left[\frac{K_h(X_i - X)}{f_X(X)} g(X) \right] \mathcal{H}\{Y_i, \theta_0(X_i, \beta_0)\} + o_P(n^{-1/2}) \\
&= n^{-1} \sum_{i=1}^n E_X \left(\frac{K_h(X_i - X)}{f_X(X)} g(X) [\mathcal{H}\{Y_i, \theta_0(X_i, \beta_0)\} - \mathcal{H}\{Y_i, \hat{\theta}(X, \beta_0)\}] \right) + o_P(n^{-1/2}),
\end{aligned}$$

since $n^{-1} \sum_{i=1}^n K_h(X_i - x) \mathcal{H}\{Y_i, \hat{\theta}(x, \beta_0)\} = o_P(n^{-1/2})$ uniformly in x , see assumption (BF5)(iii). Note that throughout this proof all expectations are conditional on the data (X_i, Y_i) , which implies that $\hat{\theta}$ is considered as constant.

Using Lemma A.1 the latter expression can be written as

$$\begin{aligned} & E_{X_1, X_2, Y_1} \left(\frac{K_h(X_1 - X_2)}{f_X(X_2)} g(X_2) [\mathcal{H}\{Y_1, \theta_0(X_1, \beta_0)\} - \mathcal{H}\{Y_1, \hat{\theta}(X_2, \beta_0)\}] \right) + o_P(n^{-1/2}) \\ &= E_{X_1, X_2} \left(\frac{K_h(X_1 - X_2)}{f_X(X_2)} g(X_2) [k\{X_1, \theta_0(X_1, \beta_0)\} - k\{X_1, \hat{\theta}(X_2, \beta_0)\}] \right) + o_P(n^{-1/2}), \end{aligned}$$

where $k(X, \theta) = E[\mathcal{H}(Y, \theta)|X]$. Using a Taylor expansion of order two and assumptions (BF1), (BF3), (BF4) and (BF6)(iv) this can be written as

$$\begin{aligned} & E_{X_2} \left(\frac{E_{X_1}\{K_h(X_1 - X_2)\}}{f_X(X_2)} g(X_2) [k\{X_2, \theta_0(X_2, \beta_0)\} - k\{X_2, \hat{\theta}(X_2, \beta_0)\}] \right) \\ &+ E_{X_2} \left(\frac{E_{X_1}\{(X_1 - X_2)K_h(X_1 - X_2)\}}{f_X(X_2)} g(X_2) \right. \\ &\quad \left. \times \frac{d}{dx} [k\{x, \theta_0(x, \beta_0)\} - k\{x, \hat{\theta}(X_2, \beta_0)\}]_{x=X_2} \right) + o_P(n^{-1/2}) \\ &= E(g(X) [k\{X, \theta_0(X, \beta_0)\} - k\{X, \hat{\theta}(X, \beta_0)\}]) + o_P(n^{-1/2}) \\ &= -E \left[g(X) \frac{\partial}{\partial \theta} k\{X, \theta_0(X, \beta_0)\} \{ \hat{\theta}(X, \beta_0) - \theta_0(X, \beta_0) \} \right] + o_P(n^{-1/2}), \end{aligned}$$

since $\sup_x |\hat{\theta}(x, \beta_0) - \theta_0(x, \beta_0)| = o_P(n^{-1/4})$. The latter expression equals $\Gamma_{BF, \theta}(\beta_0, \theta_0) [\hat{\theta} - \theta_0] + o_P(n^{-1/2})$, by using (18). Hence, the result follows.

Proof of Theorem 2.1. We will make use of Theorem 2 in Chen, Linton and Van Keilegom (2003) (CLV hereafter), which states primitive conditions under which $\hat{\beta}_{BF}$ is asymptotically normal. First of all, we need to show that $\hat{\beta}_{BF} - \beta_0 = o_P(1)$. For this, we verify the conditions of Theorem 1 in CLV. Condition (1.1) holds by definition of $\hat{\beta}_{BF}$, while the second, third and fourth condition are guaranteed by assumptions (BF5)(i) and (BF8). Finally, condition (1.5) is weaker than condition (2.5) of Theorem 2 of CLV, which we will verify below. So, the conditions of Theorem 1 are verified, up to condition (1.5) which we postpone to later. Next, we verify conditions (2.1)–(2.6) of Theorem 2 in CLV. Condition (2.1) is, as for condition (1.1), valid by construction of the estimator $\hat{\beta}_{BF}$, while condition (2.2) follows from assumption (BF6)(v). Since $\Gamma_{BF, \theta}(\beta, \theta_0) [\theta - \theta_0] = E \left\{ \frac{\partial}{\partial \theta} d(X, \theta_0) (\theta - \theta_0) (X, \beta) \right\}$, where $d(X, \theta) = E[\mathcal{L}_\beta\{Y, \beta, \theta(X, \beta)\}|X]$, we have

$$\begin{aligned} & M_{BF}(\beta, \theta) - M_{BF}(\beta, \theta_0) - \Gamma_{BF, \theta}(\beta, \theta_0) [\theta - \theta_0] \\ &= E \left\{ d(X, \theta) - d(X, \theta_0) - \frac{\partial}{\partial \theta} d(X, \theta_0) (\theta - \theta_0) (X, \beta) \right\} \end{aligned} \tag{19}$$

$$= \frac{1}{2} E \left\{ \frac{\partial^2}{\partial \theta^2} d(X, \xi) (\theta - \theta_0)^2 (X, \beta) \right\},$$

where $\xi(X)$ is in between $\theta(X, \beta)$ and $\theta_0(X, \beta)$. Hence the norm of (19) is bounded by a constant times $\|\theta - \theta_0\|_\infty^2$. This shows the first part of condition (2.3). For the second part, it follows from the proof of Theorem 2 in CLV that it suffices to show that

$$\|\Gamma_{BF, \theta}(\widehat{\beta}, \theta_0)[\widehat{\theta} - \theta_0] - \Gamma_{BF, \theta}(\beta_0, \theta_0)[\widehat{\theta} - \theta_0]\| = o_P(1) \|\widehat{\beta} - \beta_0\|,$$

and this in turn follows from (BF4), (BF5)(ii) and (BF6)(iv). Next, condition (2.4) follows from assumption (BF5)(i), while condition (2.5) is guaranteed by Theorem 3 in CLV together with assumption (BF7). It remains to verify condition (2.6). Since $\Gamma_{BF, \theta}(\beta_0, \theta_0)[\widehat{\theta} - \theta_0]$ and $M_{nBF}(\beta_0, \theta_0)$ are a sum of iid terms plus negligible terms of lower order (see Lemma A.2), this follows immediately. The asymptotic normality of $\widehat{\beta}_{BF}$ now follows.

Appendix B: Proofs for Profiling

Similarly as for the backfitting estimator, define for any $\theta \in \Theta$ and $\eta \in \Theta^q$,

$$\begin{aligned} M_{nPR}(\beta, \theta, \eta) &= n^{-1} \sum_{i=1}^n m_{PR}\{Y_i, \beta, \theta(X_i, \beta), \eta(X_i, \beta)\}, \\ M_{PR}(\beta, \theta, \eta) &= E[m_{PR}\{Y, \beta, \theta(X, \beta), \eta(X, \beta)\}], \end{aligned}$$

and let $\Gamma_{PR, \beta}(\beta, \theta, \eta) = \frac{d}{d\beta} M_{PR}(\beta, \theta, \eta)$. Note that $M_{PR}(\beta_0, \theta_0, \theta_{0\beta}) = 0$ and that

$$\begin{aligned} \Gamma_{PR, \beta}(\beta, \theta_0, \eta) &= \frac{d}{d\beta} E[\mathcal{L}_\beta\{Y, \beta, \theta_0(X, \beta)\}] + \frac{d}{d\beta} E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\} \eta(X, \beta)] \\ &= \frac{d}{d\beta} E[\mathcal{L}_\beta\{Y, \beta, \theta_0(X, \beta)\}], \end{aligned}$$

since $E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\} | X] = 0$. For functions $\xi(\cdot)$ and $\zeta(\cdot)$, let

$$\Gamma_{PR, \theta, \eta}(\beta, \theta, \eta)[\xi, \zeta] = \lim_{\tau \rightarrow 0} \frac{1}{\tau} \{M_{PR}(\beta, \theta + \tau\xi, \eta + \tau\zeta) - M_{PR}(\beta, \theta, \eta)\}.$$

Recall that Θ is some space of functions $\theta = \theta(x, \beta)$ ($x \in \mathbb{R}, \beta \in \mathcal{B}$) for which $\|\theta\|_\infty \leq M$ for some $M > 0$. For any $r \geq 1$ and any $\theta_1, \dots, \theta_r \in \Theta$, let $\|(\theta_1, \dots, \theta_r)\|_\infty = \max_{1 \leq j \leq r} \|\theta_j\|_\infty$.

The assumptions we need to impose for the main result, are the followings :

(PR1) $\theta_0 \in \Theta$, θ_0 is partially differentiable with respect to the components of β , $\frac{\partial \theta_0}{\partial \beta} \in \Theta^q$, $P(\widehat{\theta} \in \Theta) \rightarrow 1$ and $P(\widehat{\theta}_\beta \in \Theta^q) \rightarrow 1$ as $n \rightarrow \infty$, $\|\widehat{\theta} - \theta_0\|_\infty = o_P(n^{-1/4})$ and $\|\widehat{\theta}_\beta - \theta_{0\beta}\|_\infty = o_P(n^{-1/4})$.

(PR2) (i) For all y , $\mathcal{L}(y, \beta, \theta)$ is differentiable with respect to β and θ , for almost all β and θ .

(ii) $\frac{\partial}{\partial \beta} E[\mathcal{L}_\beta\{Y, \beta, \theta_0(X, \beta)\}|X]$ and $\frac{\partial}{\partial \beta} E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}|X]$ exist for all $\beta \in \mathcal{B}$, and they are equal.

(iii) $\frac{\partial^2}{\partial \theta^2} E\{\mathcal{L}_\beta(Y, \beta, \theta)|X = x\}$ and $\frac{\partial^2}{\partial \theta^2} E\{\mathcal{L}_\theta(Y, \beta, \theta)|X = x\}$ exist for all β, θ and x , and

$$\sup_{\beta \in \mathcal{B}, |\theta| \leq M, x \in R_X} \left| \frac{\partial^2}{\partial \theta^2} E\{\mathcal{L}_\beta(Y, \beta, \theta)|X = x\} \right| < \infty,$$

$$\sup_{\beta \in \mathcal{B}, |\theta| \leq M, x \in R_X} \left| \frac{\partial^2}{\partial \theta^2} E\{\mathcal{L}_\theta(Y, \beta, \theta)|X = x\} \right| < \infty,$$

where R_X is the support of X .

(iv) $\mathcal{G}(\beta)$ exists for β in a neighborhood of β_0 , is continuous at β_0 and $\mathcal{G}(\beta_0)$ is of full rank.

(PR3) (i)

$$E \left\{ \sup_{(\beta', \theta') : \|\beta - \beta'\| \leq \delta, \|\theta - \theta'\|_\infty \leq \delta, \|\eta - \eta'\|_\infty \leq \delta} |m_{PR, \ell}(Y, \beta, \theta, \eta) - m_{PR, \ell}(Y, \beta', \theta', \eta')|^{r_\ell} \right\} \leq K_\ell \delta^{r_\ell s_\ell},$$

for some $r_\ell \geq 2$, for all $(\beta, \theta, \eta) \in \mathcal{B} \times \Theta^{q+1}$, all $\delta > 0$, for some $0 < s_\ell \leq 1$ and some $K_\ell > 0$ ($\ell = 1, \dots, q$).

(ii) $\int_0^\infty \sqrt{\log N(\varepsilon^{1/s_\ell}, \widetilde{\Theta}, \|\cdot\|_\infty)} d\varepsilon < \infty$, for $\ell = 1, \dots, q$, where $\widetilde{\Theta} = \{\theta(\cdot, \beta) : \theta \in \Theta, \beta \in \mathcal{B}\}$.

(PR4)(i) For all $\delta > 0$, there exists a $\varepsilon > 0$ such that $\inf_{\|\beta - \beta_0\| > \delta} \|M_{PR}(\beta, \theta_0, \theta_{0\beta})\| \geq \varepsilon$.

(ii) Uniformly for all $\beta \in \mathcal{B}$, $M_{PR}(\beta, \theta, \eta)$ is continuous in (θ, η) at $(\theta_0, \theta_{0\beta})$ (with respect to the $\|\cdot\|_\infty$ norm).

Lemma B.1 Assume (PR1)–(PR4). Then, for any $\xi \in \Theta$, $\zeta \in \Theta^q$ and $\beta \in \mathcal{B}$,

$$\Gamma_{PR, \theta, \eta}(\beta, \theta_0, \theta_{0\beta})[\xi, \zeta] = 0.$$

Proof. Write

$$\Gamma_{PR, \theta, \eta}(\beta, \theta_0, \theta_{0\beta})[\xi, \zeta]$$

$$\begin{aligned}
&= \lim_{\tau \rightarrow 0} \frac{1}{\tau} E[\mathcal{L}_\beta\{Y, \beta, (\theta_0 + \tau\xi)(X, \beta)\} - \mathcal{L}_\beta\{Y, \beta, \theta_0(X, \beta)\}] \\
&\quad + \lim_{\tau \rightarrow 0} \frac{1}{\tau} E\{[\mathcal{L}_\theta\{Y, \beta, (\theta_0 + \tau\xi)(X, \beta)\} - \mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}](\theta_{0\beta} + \tau\zeta)(X, \beta)\} \\
&\quad + \lim_{\tau \rightarrow 0} \frac{1}{\tau} E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}\tau\zeta(X, \beta)]. \tag{20}
\end{aligned}$$

The third term of (20) equals

$$E\left(E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}|X]\zeta(X, \beta)\right) = 0,$$

since $E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}|X] = 0$. The first term of (20) can be written as

$$E\left\{\left(\frac{\partial}{\partial\theta}E[\mathcal{L}_\beta\{Y, \beta, \theta_0(X, \beta)\}|X]\right)\xi(X, \beta)\right\},$$

while the second term equals

$$E\left\{\left(\frac{\partial}{\partial\theta}E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}|X]\right)\xi(X, \beta)\frac{\partial}{\partial\beta}\theta_0(X, \beta)\right\}. \tag{21}$$

Since $E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}|X] = 0$ for all β , it follows that

$$\frac{\partial}{\partial\beta}E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}|X] + \frac{\partial}{\partial\theta}E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}|X]\frac{\partial}{\partial\beta}\theta_0(X, \beta) = 0,$$

and hence, plugging in this expression into (21) gives

$$-E\left\{\frac{\partial}{\partial\beta}E[\mathcal{L}_\theta\{Y, \beta, \theta_0(X, \beta)\}|X]\xi(X, \beta)\right\}.$$

Hence, $\Gamma_{PR, \theta, \eta}(\beta, \theta_0, \theta_{0\beta})[\xi, \zeta] = 0$, since $\frac{\partial}{\partial\beta}E(\mathcal{L}_\theta) = \frac{\partial}{\partial\theta}E(\mathcal{L}_\beta)$.

Proof of Theorem 2.2. In a manner similar to the backfitting procedure, we proceed by checking the primitive conditions of Theorem 2 in Chen, Linton and Van Keilegom (2003) (CLV hereafter). Note that the results in that paper are valid for one-dimensional nuisance functions θ , but it is readily seen how to extend their primitive conditions to the current setup of $(q + 1)$ -dimensional nuisance functions.

The verification of the conditions in that theorem is much the same as for the backfitting procedure, except for conditions (2.3) and (2.5). Let us start with verifying (2.3). Since it follows from the proof of Lemma B.1 that $\Gamma_{PR, \theta, \eta}(\beta, \theta_0, \theta_{0\beta})[\theta - \theta_0, \eta - \theta_{0\beta}] =$

$E\{\frac{\partial}{\partial\theta}d_1(X, \theta_0)(\theta - \theta_0)(X, \beta)\} + E\{\frac{\partial}{\partial\theta}d_2(X, \theta_0)(\theta - \theta_0)(X, \beta)\frac{\partial}{\partial\beta}\theta_0(X, \beta)\}$, where $d_1(X, \theta) = E[\mathcal{L}_\beta\{Y, \beta, \theta(X, \beta)\}|X]$ and $d_2(X, \theta) = E[\mathcal{L}_\theta\{Y, \beta, \theta(X, \beta)\}|X]$, we have

$$\begin{aligned}
& M_{PR}(\beta, \theta, \eta) - M_{PR}(\beta, \theta_0, \theta_{0\beta}) - \Gamma_{PR, \theta, \eta}(\beta, \theta_0, \theta_{0\beta})[\theta - \theta_0, \eta - \theta_{0\beta}] \tag{22} \\
&= E\left\{d_1(X, \theta) - d_1(X, \theta_0) - \frac{\partial}{\partial\theta}d_1(X, \theta_0)(\theta - \theta_0)(X, \beta)\right\} \\
&+ E\left[\left\{d_2(X, \theta) - d_2(X, \theta_0) - \frac{\partial}{\partial\theta}d_2(X, \theta_0)(\theta - \theta_0)(X, \beta)\right\}\eta(X, \beta)\right] \\
&+ E\left\{d_2(X, \theta_0)(\eta - \theta_{0\beta})(X, \beta)\right\} \\
&+ E\left\{\frac{\partial}{\partial\theta}d_2(X, \theta_0)(\theta - \theta_0)(X, \beta)(\eta - \theta_{0\beta})(X, \beta)\right\} \\
&= \frac{1}{2}E\left\{\frac{\partial^2}{\partial\theta^2}d_1(X, \xi_1)(\theta - \theta_0)^2(X, \beta)\right\} + \frac{1}{2}E\left\{\frac{\partial^2}{\partial\theta^2}d_2(X, \xi_2)(\theta - \theta_0)^2(X, \beta)\frac{\partial}{\partial\beta}\theta_0(X, \beta)\right\} \\
&+ E\left\{\frac{\partial}{\partial\theta}d_2(X, \theta_0)(\theta - \theta_0)(X, \beta)(\eta - \theta_{0\beta})(X, \beta)\right\},
\end{aligned}$$

since $d_2(X, \theta_0) \equiv 0$, where $\xi_1(X)$ and $\xi_2(X)$ are in between $\theta(X, \beta)$ and $\theta_0(X, \beta)$. Hence the norm of (22) is bounded by a constant times $\|(\theta - \theta_0, \eta - \theta_{0\beta})\|_\infty^2$. This shows the first part of condition (2.3). The second part is obvious by Lemma B.1.

Finally, condition (2.5) is guaranteed by Theorem 3 in CLV together with assumption (PR3). Note that $N(\varepsilon^{1/s_\varepsilon}, \tilde{\Theta}^q, \|\cdot\|_\infty) \leq N(\varepsilon^{1/s_\varepsilon}, \tilde{\Theta}, \|\cdot\|_\infty)^q$ and hence the second condition in Theorem 3 in CLV is implied by (PR3)(ii). The result now follows.

References

- Akritis, M.G. and Van Keilegom, I. (2001). Nonparametric estimation of the residual distribution. *Scand. J. Statist.*, **28**, 549–568.
- Buja, A., Hastie, T. J. and Tibshirani, R. J. (1989). Linear smoothers and additive models (with Discussion). *Annals of Statistics*, **17**, 453–555.
- Carroll, R. J., Fan, J., Gijbels, I. and Wand, M. P. (1997). Generalized partially linear single-index models. *Journal of the American Statistical Association*, **92**, 477–489.
- Chaudhuri, P. (1991). Nonparametric estimates of regression quantiles and their local Bahadur representation. *Annals of Statistics*, **19**, 760–777.

- Chen, X., Linton, O. and Van Keilegom, I. (2003). Estimation of semiparametric models when the criterion function is not smooth. *Econometrica*, **71**, 1591–1608.
- Feder, P. I. (1975). On asymptotic distribution theory in segmented regression problems - identified case. *Annals of Statistics*, **3**, 49–83.
- Hastie, T., and Tibshirani, R. J. (1990). *Generalized Additive Models*. London: Chapman and Hall.
- Hu, Z., Wang, N. and Carroll, R. J. (2004). Profile-kernel versus backfitting in the partially linear model for longitudinal/clustered data. *Biometrika*, **91**, 251–262.
- Mammen, E., Linton, O. and Nielsen, J. (1999). The existence and asymptotic properties of a backfitting projection algorithm under weak conditions. *Annals of Statistics*, **27**, 1443–1490.
- Murphy, S. A. and van der Vaart, A. W. (2000). On profile likelihood. *Journal of the American Statistical Association*, **95**, 449–485.
- Opsomer, J. D. (2000). Asymptotic properties of backfitting estimators. *Journal of Multivariate Analysis*, **73**, 166–179.
- Opsomer, J. D. and Ruppert, D. (1997). Fitting a bivariate additive model by local polynomial regression. *Annals of Statistics*, **25**, 186–211.
- Opsomer, J. D. and Ruppert, D. (1999). A root-n consistent backfitting estimator for semi-parametric additive modeling. *Journal of Computational and Graphical Statistics*, **8**, 715–732.
- Rice, J. A. (1986). Convergence rates for partially splined models. *Statistics and Probability Letters*, **4**, 204–208.
- Roberts, L. (1991). Dioxin risks revisited. *Science*, **251**, 624–626.
- Severini, T. A. and Staniswalis, J. G. (1994). Quasi-likelihood estimation in semiparametric models. *Journal of the American Statistical Association*, **89**, 501–512.
- Severini, T. A. & Wong, W. H. (1992). Profile likelihood and conditionally parametric

- models. *Annals of Statistics*, **20**, 1768–1802.
- Speckman, P. E. (1988). Regression analysis for partially linear models. *Journal of the Royal Statistical Society, Series B*, **50**, 413–436.
- van der Vaart, A. W. and Wellner, J. A. (1996). *Weak Convergence and Empirical Processes*. Springer-Verlag, New York.
- Wand, M. P. (1999). A central limit theorem for local polynomial backfitting estimators. *Journal of Multivariate Analysis*, **70**, 57–65.