

Robust-BD Estimation and Inference for General Partially Linear Models

Chunming Zhang * and Zhengjun Zhang

Department of Statistics, University of Wisconsin-Madison, Madison, WI 53706, USA; zjz@stat.wisc.edu

* Correspondence: cmzhang@stat.wisc.edu

Received: 10 October 2017; Accepted: 16 November 2017; Published: 20 November 2017

Abstract: The classical quadratic loss for the partially linear model (PLM) and the likelihood function for the generalized PLM are not resistant to outliers. This inspires us to propose a class of “robust-Bregman divergence (BD)” estimators of both the parametric and nonparametric components in the general partially linear model (GPLM), which allows the distribution of the response variable to be partially specified, without being fully known. Using the local-polynomial function estimation method, we propose a computationally-efficient procedure for obtaining “robust-BD” estimators and establish the consistency and asymptotic normality of the “robust-BD” estimator of the parametric component β_o . For inference procedures of β_o in the GPLM, we show that the Wald-type test statistic W_n constructed from the “robust-BD” estimators is asymptotically distribution free under the null, whereas the likelihood ratio-type test statistic Λ_n is not. This provides an insight into the distinction from the asymptotic equivalence (Fan and Huang 2005) between W_n and Λ_n in the PLM constructed from profile least-squares estimators using the non-robust quadratic loss. Numerical examples illustrate the computational effectiveness of the proposed “robust-BD” estimators and robust Wald-type test in the appearance of outlying observations.

Keywords: Bregman divergence; generalized linear model; local-polynomial regression; model check; nonparametric test; quasi-likelihood; semiparametric model; Wald statistic

1. Introduction

Semiparametric models, such as the partially linear model (PLM) and generalized PLM, play an important role in statistics, biostatistics, economics and engineering studies [1–5]. For the response variable Y and covariates (X, T) , where $X = (X_1, \dots, X_d)^T \in \mathbb{R}^d$ and $T \in \mathcal{T} \subseteq \mathbb{R}^D$, the PLM, which is widely used for continuous responses Y , describes the model structure according to:

$$Y = X^T \beta_o + \eta^o(T) + \epsilon, \quad E(\epsilon | X, T) = 0, \quad (1)$$

where $\beta_o = (\beta_{1;o}, \dots, \beta_{d;o})^T$ is a vector of unknown parameters and $\eta^o(\cdot)$ is an unknown smooth function; the generalized PLM, which is more suited to discrete responses Y and extends the generalized linear model [6], assumes:

$$m(x, t) = E(Y | X = x, T = t) = F^{-1}(x^T \beta_o + \eta^o(t)), \quad (2)$$

$$Y | (X, T) \sim \text{exponential family of distributions}, \quad (3)$$

where F is a known link function. Typically, the parametric component β_o is of primary interest, while the nonparametric component $\eta^o(\cdot)$ serves as a nuisance function. For illustration clarity, this paper focuses on $D = 1$. An important application of PLM to brain fMRI data was given in [7] for detecting activated brain voxels in response to external stimuli. There, β_o corresponds to the part of hemodynamic response values, which is the object of primary interest to neuroscientists; $\eta^o(\cdot)$ is the

slowly drifting baseline of time. Determining whether a voxel is activated or not can be formulated as testing for the linear form of hypotheses,

$$H_0 : A\beta_o = g_0 \quad \text{versus} \quad H_1 : A\beta_o \neq g_0, \quad (4)$$

where A is a given $k \times d$ full row rank matrix and g_0 is a known $k \times 1$ vector.

Estimation of the parametric and nonparametric components of PLM and generalized PLM has received much attention in the literature. On the other hand, the existing work has some limitations: (i) The generalized PLM assumes that $Y | (X, T)$ follows the distribution in (3), so that the likelihood function is fully available. From the practical viewpoint, results from the generalized PLM are not applicable to situations where the distribution of $Y | (X, T)$ either departs from (3) or is incompletely known. (ii) Some commonly-used error measures, such as the quadratic loss in PLM for Gaussian-type responses (see for example [7,8]) and the (negative) likelihood function used in the generalized PLM, are not resistant to outliers. The work in [9] studied robust inference based on the kernel regression method for the generalized PLM with a canonical link, based on either the (negative) likelihood or (negative) quasi-likelihood as the error measure, and illustrated numerical examples with the dimension $d = 1$. However, the quasi-likelihood is not suitable for the exponential loss function (defined in Section 2.1), commonly used in machine learning and data mining. (iii) The work in [8] developed the inference of (4) for PLM, via the classical quadratic loss as the error measure, and demonstrated that the asymptotic distributions of the likelihood ratio-type statistic and Wald statistic under the null of (4) are both χ_k^2 . It remains unknown whether this conclusion holds when the tests are constructed based on robust estimators.

Without completely specifying the distribution of $Y | (X, T)$, we assume:

$$\text{var}(Y | X = x, T = t) = V(m(x, t)), \quad (5)$$

with a known functional form of $V(\cdot)$. We refer to a model specified by (2) and (5) as the “general partially linear model” (GPLM). This paper aims to develop robust estimation of GPLM and robust inference of β_o , allowing the distribution of $Y | (X, T)$ to be partially specified. To introduce robust estimation, we adopt a broader class of robust error measures, called “robust-Bregman divergence (BD)” developed in [10], for a GLM, in which BD includes the quadratic loss, the (negative) quasi-likelihood, the exponential loss and many other commonly-used error measures as special cases. We propose the “robust-BD estimators” for both the parametric and nonparametric components of the GPLM. Distinct from the explicit-form estimators for PLM using the classical quadratic loss (see [8]), the “robust-BD estimators” for GPLM do not have closed-form expressions, which makes the theoretical derivation challenging. Moreover, the robust-BD estimators, as numerical solutions to non-linear optimization problems, pose key implementation challenges. Our major contributions are given below.

- The robust fitting of the nonparametric component $\eta^o(\cdot)$ is formulated using the local-polynomial regression technique [11]. See Section 2.3.
- We develop a coordinate descent algorithm for the robust-BD estimator of β_o , which is computationally efficient particularly when the dimension d is large. See Section 3.
- Theorems 1 and 2 demonstrate that under the GPLM, the consistency and asymptotic normality of the proposed robust-BD estimator for β_o are achieved. See Section 4.
- For robust inference of β_o , we propose a robust version of the Wald-type test statistic W_n , based on the robust-BD estimators, and justify its validity in Theorems 3–5. It is shown to be asymptotically χ^2 (central) under the null, thus distribution free, and χ^2 (noncentral) under the contiguous alternatives. Hence, this result, when applied to the exponential loss, as well as other loss functions in the wider class of BD, is practically feasible. See Section 5.1.

- For robust inference of β_o , we re-examine the likelihood ratio-type test statistic Λ_n , constructed by replacing the negative log-likelihood with the robust-BD. Our Theorem 6 reveals that the asymptotic null distribution of Λ_n is generally not χ^2 , but a linear combination of independent χ^2 variables, with weights relying on unknown quantities. Even in the particular case of using the classical-BD, the limit distribution is not invariant with re-scaling the generating function of the BD. Moreover, the limit null distribution of Λ_n (in either the non-robust or robust version) using the exponential loss, which does not belong to the (negative) quasi-likelihood, but falls in BD, is always a weighted χ^2 , thus limiting its use in practical applications. See Section 5.2.

Simulation studies in Section 6 demonstrate that the proposed class of robust-BD estimators and robust Wald-type test either compare well with or perform better than the classical non-robust counterparts: the former is less sensitive to outliers than the latter, and both perform comparably well for non-contaminated cases. Section 7 illustrates some real data applications. Section 8 ends the paper with brief discussions. Details of technical derivations are relegated to Appendix A.

2. Robust-BD and Robust-BD Estimators

This section starts with a brief review of BD in Section 2.1 and “robust-BD” in Section 2.2, followed by the proposed “robust-BD” estimators of $\eta^o(\cdot)$ and β_o in Sections 2.3 and 2.4.

2.1. Classical-BD

To broaden the scope of robust estimation and inference, we consider a class of error measures motivated from the Bregman divergence (BD). For a given concave q -function, [12] defined a bivariate function,

$$Q_q(v, \mu) = -q(v) + q(\mu) + (v - \mu)q'(\mu). \quad (6)$$

We call Q_q the BD and call q the generating q -function of the BD. For example, a function $q(\mu) = a\mu - \mu^2$ for some constant a yields the quadratic loss $Q_q(Y, \mu) = (Y - \mu)^2$. For a binary response variable Y , $q(\mu) = \min\{\mu, (1 - \mu)\}$ gives the misclassification loss $Q_q(Y, \mu) = I\{Y \neq I(\mu > 1/2)\}$, where $I(\cdot)$ is an indicator function; $q(\mu) = -2\{\mu \log(\mu) + (1 - \mu) \log(1 - \mu)\}$ gives the Bernoulli deviance loss log-likelihood $Q_q(Y, \mu) = -2\{Y \log(\mu) + (1 - Y) \log(1 - \mu)\}$; $q(\mu) = 2 \min\{\mu, (1 - \mu)\}$ results in the hinge loss $Q_q(Y, \mu) = \max\{1 - (2Y - 1) \text{sign}(\mu - 0.5), 0\}$ of the support vector machine; $q(\mu) = 2\{\mu(1 - \mu)\}^{1/2}$ yields the exponential loss $Q_q(Y, \mu) = \exp[-(Y - 0.5) \log\{\mu/(1 - \mu)\}]$ used in AdaBoost [13]. Moreover, [14] showed that if:

$$q(\mu) = \int_a^\mu \frac{s - \mu}{V(s)} ds, \quad (7)$$

with a finite constant a such that the integral is well defined, then $Q_q(y, \mu)$ matches the “classical (negative) quasi-likelihood” function.

2.2. Robust-BD $\rho_q(y, \mu)$

Let $r(y, \mu) = (y - \mu) / \sqrt{V(\mu)}$ denote the Pearson residual, which reduces to the standardized residual for linear models. In contrast to the “classical-BD”, denoted by Q_q in (6), the “robust-BD” developed in [10] for a GLM [6], is formed by:

$$\rho_q(y, \mu) = \int_y^\mu \psi(r(y, s)) \{q''(s) \sqrt{V(s)}\} ds - G(\mu), \quad (8)$$

where $\psi(r)$ is chosen to be a bounded, odd function, such as the Huber ψ -function [15], $\psi(r) = r \min(1, c/|r|)$, and the bias-correction term, $G(\mu)$, entails the Fisher consistency of the parameter estimator and satisfies:

$$G'(\mu) = G'_1(\mu) \{q''(\mu) \sqrt{V(\mu)}\},$$

with

$$G'_1(m(\mathbf{x}, t)) = E\{\psi(r(Y, m(\mathbf{x}, t))) \mid \mathbf{X} = \mathbf{x}, T = t\}. \quad (9)$$

We make the following discussions regarding features of the “robust-BD”. To facilitate the discussion, we first introduce some necessary notation. Assume that the quantities:

$$p_j(y; \theta) = \frac{\partial^j}{\partial \theta^j} \rho_q(y, F^{-1}(\theta)), \quad j = 0, 1, \dots, \quad (10)$$

exist finitely up to any order required. Then, we have the following expressions,

$$\begin{aligned} p_1(y; \theta) &= \{\psi(r(y, \mu)) - G'_1(\mu)\} \{q''(\mu) \sqrt{V(\mu)}\} / F'(\mu), \\ p_2(y; \theta) &= A_0(y, \mu) + \{\psi(r(y, \mu)) - G'_1(\mu)\} A_1(\mu), \\ p_3(y; \theta) &= A_2(y, \mu) + \{\psi(r(y, \mu)) - G'_1(\mu)\} A'_1(\mu) / F'(\mu), \end{aligned} \quad (11)$$

where $\mu = F^{-1}(\theta)$,

$$A_0(y, \mu) = -\left[\psi'(r(y, \mu))\left\{1 + \frac{y - \mu}{\sqrt{V(\mu)}} \times \frac{V'(\mu)}{2\sqrt{V(\mu)}}\right\} + G''_1(\mu)\sqrt{V(\mu)}\right] \frac{q''(\mu)}{\{F'(\mu)\}^2},$$

$A_1(\mu) = [\{q^{(3)}(\mu)\sqrt{V(\mu)} + 2^{-1}q''(\mu)V'(\mu)/\sqrt{V(\mu)}\}F'(\mu) - q''(\mu)\sqrt{V(\mu)}F''(\mu)]/\{F'(\mu)\}^3$ and $A_2(y, \mu) = [\partial A_0(y, \mu)/\partial \mu + \partial\{\psi(r(y, \mu)) - G'_1(\mu)\}/\partial \mu A_1(\mu)]/F'(\mu)$. Particularly, $p_1(y; \theta)$ contains $\psi(r)$; $p_2(y; \theta)$ contains $\psi(r)$, $\psi'(r)$ and $\psi'(r)r$; $p_3(y; \theta)$ contains $\psi(r)$, $\psi'(r)$, $\psi'(r)r$, $\psi''(r)$, $\psi''(r)r$, and $\psi''(r)r^2$, where $r = r(y, \mu) = (y - \mu)/\sqrt{V(\mu)}$ denotes the Pearson residual. Accordingly, $\{p_j(y; \theta) : j = 1, 2, 3\}$ depend on y through $\psi(r)$ and its derivatives coupled with r . Then, we observe from (9) and (11) that:

$$E\{p_1(Y; \mathbf{X}^T \boldsymbol{\beta}_0 + \eta^0(T)) \mid \mathbf{X}, T\} = 0. \quad (12)$$

In the particular choice of $\psi(r) = r$, it is clearly noticed from (9) that $G'_1(\cdot) = 0$, and thus, $G'(\cdot) = 0$. In such a case, the proposed “robust-BD” $\rho_q(y, \mu)$ reduces to the “classical-BD” $Q_q(y, \mu)$.

2.3. Local-Polynomial Robust-BD Estimator of $\eta^0(\cdot)$

Let $\{(Y_i, X_i, T_i)\}_{i=1}^n$ be i.i.d. observations of (Y, X, T) captured by the GPLM in (2) and (5), where the dimension $d \geq 1$ is a finite integer. From (2), it is directly observed that if the true value of $\boldsymbol{\beta}_0$ is known, then estimating $\eta^0(\cdot)$ becomes estimating a nonparametric function; conversely, if the actual form of $\eta^0(\cdot)$ is available, then estimating $\boldsymbol{\beta}_0$ amounts to estimating a vector parameter.

To motivate the estimation of $\eta^0(\cdot)$ at a fitting point t , a proper way to characterize $\eta^0(t)$ is desired. For any given value of $\boldsymbol{\beta}$, define:

$$S(a; t, \boldsymbol{\beta}) = E\{\rho_q(Y, F^{-1}(\mathbf{X}^T \boldsymbol{\beta} + a)) w_1(\mathbf{X}) \mid T = t\}, \quad (13)$$

where a is a scalar, $\rho_q(y, \mu)$ is the “robust-BD” defined in (8), which aims to guard against outlying observations in the response space of Y , and $w_1(\cdot) \geq 0$ is a given bounded weight function that downweights high leverage points in the covariate space of X . See Sections 6 and 7 for an example of $w_1(\mathbf{x})$. Set:

$$\eta_{\boldsymbol{\beta}}(t) = \arg \min_{a \in \mathbb{R}^1} S(a; t, \boldsymbol{\beta}). \quad (14)$$

Theoretically, $\eta^0(t) = \eta_{\boldsymbol{\beta}_0}(t)$ will be assumed (in Condition A3) for obtaining asymptotically unbiased estimators of $\eta^0(\cdot)$. Such property indeed holds, for example, when a classical quadratic loss combined with an identity link is used in (14). Thus, we call $\eta_{\boldsymbol{\beta}}(\cdot)$ the “surrogate function” for $\eta^0(\cdot)$.

The characterization of the surrogate function $\eta_{\boldsymbol{\beta}}(t)$ in (14) enables us to develop its robust-BD estimator $\hat{\eta}_{\boldsymbol{\beta}}(t)$ based on nonparametric function estimation. Assume that

$\eta^0(\cdot)$ is $(p+1)$ -times continuously differentiable at the fitting point t . Denote by $\mathbf{a}_0(t) = (\eta^0(t), (\eta^0)^{(1)}(t), \dots, (\eta^0)^{(p)}(t)/p!)^T \in \mathbb{R}^{p+1}$ the vector consisting of $\eta^0(t)$ along with its (re-scaled) derivatives. For observed covariates T_i close to the point t , the Taylor expansion implies that:

$$\begin{aligned}\eta^0(T_i) &\approx \eta^0(t) + (T_i - t)(\eta^0)^{(1)}(t) + \dots + (T_i - t)^p(\eta^0)^{(p)}(t)/p! \\ &= \mathbf{t}_i(t)^T \mathbf{a}_0(t),\end{aligned}\quad (15)$$

where $\mathbf{t}_i(t) = (1, (T_i - t), \dots, (T_i - t)^p)^T$. For any given value of $\boldsymbol{\beta}$, let $\hat{\mathbf{a}}(t; \boldsymbol{\beta}) = (\hat{a}_0(t; \boldsymbol{\beta}), \hat{a}_1(t; \boldsymbol{\beta}), \dots, \hat{a}_p(t; \boldsymbol{\beta}))^T$ be the minimizer of the criterion function,

$$S_n(\mathbf{a}; t, \boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \rho_q(Y_i, F^{-1}(\mathbf{X}_i^T \boldsymbol{\beta} + \mathbf{t}_i(t)^T \mathbf{a})) w_1(\mathbf{X}_i) K_h(T_i - t), \quad (16)$$

with respect to $\mathbf{a} \in \mathbb{R}^{p+1}$, where $K_h(\cdot) = K(\cdot/h)/h$ is re-scaled from a kernel function K and $h > 0$ is termed a bandwidth parameter. The first entry of $\hat{\mathbf{a}}(t; \boldsymbol{\beta})$ supplies the local-polynomial robust-BD estimator $\hat{\eta}_{\boldsymbol{\beta}}(t)$ of $\eta_{\boldsymbol{\beta}}(t)$, i.e.,

$$\hat{\eta}_{\boldsymbol{\beta}}(t) = \mathbf{e}_{1,p+1}^T \left\{ \arg \min_{\mathbf{a} \in \mathbb{R}^{p+1}} S_n(\mathbf{a}; t, \boldsymbol{\beta}) \right\}, \quad (17)$$

where $\mathbf{e}_{j,p+1}$ denotes the j -th column of a $(p+1) \times (p+1)$ identity matrix.

It is noted that the reliance of $\hat{\eta}_{\boldsymbol{\beta}}(t)$ on $\boldsymbol{\beta}$ does not guarantee its consistency to $\eta^0(t)$. Nonetheless, it is anticipated from the uniform consistency of $\hat{\eta}_{\hat{\boldsymbol{\beta}}}$ in Lemma 1 that $\hat{\eta}_{\hat{\boldsymbol{\beta}}}(t)$ will offer a valid estimator of $\eta^0(t)$, provided that $\hat{\boldsymbol{\beta}}$ consistently estimates $\boldsymbol{\beta}_0$. Section 2.4 will discuss our proposed robust-BD estimator $\hat{\boldsymbol{\beta}}$. Furthermore, Lemma 1 will assume (in Condition A1) that $\eta_{\boldsymbol{\beta}}(t)$ is the unique minimizer of $S(\mathbf{a}; t, \boldsymbol{\beta})$ with respect to \mathbf{a} .

Remark 1. The case of using the “kernel estimation”, or locally-constant estimation, corresponds to the choice of degree $p = 0$ in (15). In that case, the criterion function in (16) and the estimator in (17) reduce to:

$$S_n(\mathbf{a}; t, \boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \rho_q(Y_i, F^{-1}(\mathbf{X}_i^T \boldsymbol{\beta} + \mathbf{a})) w_1(\mathbf{X}_i) K_h(T_i - t), \quad (18)$$

$$\hat{\eta}_{\boldsymbol{\beta}}(t) = \arg \min_{\mathbf{a} \in \mathbb{R}^1} S_n(\mathbf{a}; t, \boldsymbol{\beta}), \quad (19)$$

respectively.

2.4. Robust-BD Estimator of $\boldsymbol{\beta}_0$

For any given value of $\boldsymbol{\beta}$, define:

$$J(\boldsymbol{\beta}, \eta_{\boldsymbol{\beta}}) = E\{\rho_q(Y, F^{-1}(\mathbf{X}^T \boldsymbol{\beta} + \eta_{\boldsymbol{\beta}}(T))) w_2(\mathbf{X})\}, \quad (20)$$

where $\eta_{\boldsymbol{\beta}}(\cdot)$ is as defined in (14) and $w_2(\cdot)$ plays the same role as $w_1(\cdot)$ in (13). Theoretically, it is anticipated that:

$$\boldsymbol{\beta}_0 = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^d} J(\boldsymbol{\beta}, \eta_{\boldsymbol{\beta}}), \quad (21)$$

which holds for example in the case where a classical quadratic loss combined with an identity link is used. To estimate $\boldsymbol{\beta}_0$, it is natural to replace (20) by its sample-based criterion,

$$J_n(\boldsymbol{\beta}, \hat{\eta}_{\boldsymbol{\beta}}) = \frac{1}{n} \sum_{i=1}^n \rho_q(Y_i, F^{-1}(\mathbf{X}_i^T \boldsymbol{\beta} + \hat{\eta}_{\boldsymbol{\beta}}(T_i))) w_2(\mathbf{X}_i), \quad (22)$$

where $\hat{\eta}_\beta(\cdot)$ is as defined in (17). Hence, a parametric estimator of β_0 is provided by:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} J_n(\beta, \hat{\eta}_\beta). \quad (23)$$

Finally, the estimator of $\eta^0(\cdot)$ is given by:

$$\hat{\eta}(\cdot) = \hat{\eta}_{\hat{\beta}}(\cdot).$$

To achieve asymptotic normality of $\hat{\beta}$, Theorem 2 assumes (in Condition A2) that β_0 is the unique minimizer in (21), a standard condition for consistent M -estimators [16].

As a comparison, it is seen that $w_1(\cdot)$ in (16) is used to robustify covariates X_i in estimating $\eta^0(\cdot)$, $w_2(\cdot)$ in (22) is used to robustify covariates X_i in estimating β_0 and $\rho_q(\cdot, \cdot)$ serves to robustify the responses Y_i in both estimating procedures.

3. Two-Step Iterative Algorithm for Robust-BD Estimation

In a special case of using the classical quadratic loss combined with an identity link function, the robust-BD estimators for parametric and nonparametric components have explicit expressions,

$$\hat{\beta} = (\tilde{X}^T \mathbf{w}_2 \tilde{X})^{-1} (\tilde{X}^T \mathbf{w}_2 \tilde{y}), \quad (\hat{\eta}(T_1), \dots, \hat{\eta}(T_n))^T = S_h(\mathbf{y} - \mathbf{X}\hat{\beta}), \quad (24)$$

where $\mathbf{w}_2 = \text{diag}(w_2(X_1), \dots, w_2(X_n))$, $\tilde{y} = (\mathbf{I} - S_h)\mathbf{y}$, $\tilde{X} = (\mathbf{I} - S_h)\mathbf{X}$, with \mathbf{I} being an identity matrix, $\mathbf{y} = (Y_1, \dots, Y_n)^T$, $\mathbf{X} = (X_1, \dots, X_n)^T$ the design matrix,

$$S_h = \begin{pmatrix} e_{1,p+1}^T [\{\mathbf{T}(T_1)\}^T \mathbf{W}_{w_1;K}(T_1) \mathbf{T}(T_1)]^{-1} \{\mathbf{T}(T_1)\}^T \mathbf{W}_{w_1;K}(T_1) \\ \vdots \\ e_{1,p+1}^T [\{\mathbf{T}(T_n)\}^T \mathbf{W}_{w_1;K}(T_n) \mathbf{T}(T_n)]^{-1} \{\mathbf{T}(T_n)\}^T \mathbf{W}_{w_1;K}(T_n) \end{pmatrix},$$

and:

$$\mathbf{T}(t) = (\mathbf{t}_1(t), \dots, \mathbf{t}_n(t))^T, \quad \mathbf{W}_{w_1;K}(t) = \text{diag}\{w_1(X_i) K_h(T_i - t) : i = 1, \dots, n\}.$$

When $w_1(x) = w_2(x) \equiv 1$, (24) reduces to the “profile least-squares estimators” of [8].

In other cases, robust-BD estimators from (17) and (23) do not have closed-form expressions and need to be solved numerically, which are computationally challenging and intensive. We now discuss a two-step robust proposal for iteratively estimating β_0 and $\eta^0(\cdot)$. Let $\hat{\beta}^{[k-1]}$ and $\{\hat{\eta}^{[k-1]}(T_i)\}_{i=1}^n$ denote the estimates in the $(k-1)$ -th iteration, where $\hat{\eta}^{[k-1]}(\cdot) = \hat{\eta}_{\hat{\beta}^{[k-1]}}(\cdot)$. The k -th iteration consists of two steps below.

Step 1: Instead of solving (23) directly, we propose to solve a surrogate optimization problem, $\hat{\beta}^{[k]} = \arg \min_{\beta \in \mathbb{R}^d} J_n(\beta, \hat{\eta}^{[k-1]})$. This minimizer approximates $\hat{\beta}$.

Step 2: Obtain $\hat{\eta}^{[k]}(T_i) = \hat{\eta}_{\hat{\beta}^{[k]}}(T_i)$, $i = 1, \dots, n$, where $\hat{\eta}_\beta(t)$ is defined in (17).

The algorithm terminates provided that $\|\hat{\beta}^{[k]} - \hat{\beta}^{[k-1]}\|$ is below some pre-specified threshold value, and all $\{\hat{\eta}^{[k]}(T_i)\}_{i=1}^n$ stabilize.

3.1. Step 1

For the above two-step algorithm, we first elaborate on the procedure of acquiring $\hat{\beta}^{[k]}$ in Step 1, by extending the coordinate descent (CD) iterative algorithm [17] designed for penalized estimation to our current robust-BD estimation, which is computationally efficient. For any given value of η ,

by Taylor expansion, around some initial estimate β^* (for example, $\hat{\beta}^{[k-1]}$), we obtain the weighted quadratic approximation,

$$\rho_q(Y_i, F^{-1}(X_i^T \beta + \eta)) \approx \frac{1}{2} s_i^I (Z_i^I - X_i^T \beta)^2 + C_i,$$

where C_i is a constant not depending on β ,

$$\begin{aligned} s_i^I &= p_2(Y_i; X_i^T \beta^* + \eta), \\ Z_i^I &= X_i^T \beta^* - p_1(Y_i; X_i^T \beta^* + \eta) / p_2(Y_i; X_i^T \beta^* + \eta), \end{aligned}$$

with $p_j(y; \theta)$ defined in (10). Hence,

$$\begin{aligned} J_n(\beta, \eta) &= \frac{1}{n} \sum_{i=1}^n \rho_q(Y_i, F^{-1}(X_i^T \beta + \eta)) w_2(X_i) \\ &\approx \frac{1}{2} \sum_{i=1}^n \left\{ n^{-1} s_i^I w_2(X_i) \right\} (Z_i^I - X_i^T \beta)^2 + \text{constant}. \end{aligned}$$

Thus it suffices to conduct minimization of $\sum_{i=1}^n s_i^I w_2(X_i) (Z_i^I - X_i^T \beta)^2$ with respect to β , using a coordinate descent (CD) updating procedure. Suppose that the current estimate is $\hat{\beta}^{\text{old}} = (\hat{\beta}_1^{\text{old}}, \dots, \hat{\beta}_d^{\text{old}})^T$, with the current residual vector $\hat{r}^{\text{old}} = (\hat{r}_1^{\text{old}}, \dots, \hat{r}_n^{\text{old}})^T = \mathbf{z}^I - \mathbf{X} \hat{\beta}^{\text{old}}$, where $\mathbf{z}^I = (Z_1^I, \dots, Z_n^I)^T$ is the vector of pseudo responses. Adopting the Newton–Raphson algorithm, the estimate of the j -th coordinate based on the previous estimate $\hat{\beta}_j^{\text{old}}$ is updated to:

$$\hat{\beta}_j^{\text{new}} = \hat{\beta}_j^{\text{old}} + \frac{\sum_{i=1}^n \{s_i^I w_2(X_i)\} \hat{r}_i^{\text{old}} X_{i,j}}{\sum_{i=1}^n \{s_i^I w_2(X_i)\} X_{i,j}^2}.$$

As a result, the residuals due to such an update are updated to:

$$\hat{r}_i^{\text{new}} = \hat{r}_i^{\text{old}} - X_{i,j} (\hat{\beta}_j^{\text{new}} - \hat{\beta}_j^{\text{old}}), \quad i = 1, \dots, n.$$

Cycling through $j = 1, \dots, d$, we obtain the estimate $\hat{\beta}^{\text{new}} = (\hat{\beta}_1^{\text{new}}, \dots, \hat{\beta}_d^{\text{new}})^T$. Now, we set $\eta = \hat{\eta}^{[k-1]}$ and $\beta^* = \hat{\beta}^{[k-1]}$. Iterate the process of weighted quadratic approximation followed by the CD updating, for a number of times, until the estimate $\hat{\beta}^{\text{new}}$ stabilizes to the solution $\hat{\beta}^{[k]}$.

The validity of $\hat{\beta}^{[k]}$ in Step 1 converging to the true parameter β_o is justified as follows. (i) Standard results for M -estimation [16] indicate that the minimizer of $J_n(\beta, \eta_{\beta_o})$ is consistent with β_o . (ii) According to our Theorem 1 (ii) in Section 4.1, $\sup_{t \in \mathcal{T}} |\hat{\eta}_{\hat{\beta}}(t) - \eta_{\beta_o}(t)| \xrightarrow{P} 0$ for a compact set \mathcal{T} , where \xrightarrow{P} stands for convergence in probability. Using derivations similar to those of (A4) gives $\sup_{\beta \in \mathcal{K}} |J_n(\beta, \hat{\eta}_{\hat{\beta}}) - J_n(\beta, \eta_{\beta_o})| \xrightarrow{P} 0$ for any compact set \mathcal{K} . Thus, minimizing $J_n(\beta, \hat{\eta}_{\hat{\beta}})$ is asymptotically equivalent to minimizing $J_n(\beta, \eta_{\beta_o})$. (iii) Similarly, provided that $\hat{\beta}^{[k-1]}$ is close to $\hat{\beta}$, minimizing $J_n(\beta, \hat{\eta}_{\hat{\beta}^{[k-1]}})$ is asymptotically equivalent to minimizing $J_n(\beta, \hat{\eta}_{\hat{\beta}})$. Assembling these three results with the definition of $\hat{\beta}^{[k]}$ yields:

$$\begin{aligned} \hat{\beta}^{[k]} &= \arg \min_{\beta} J_n(\beta, \hat{\eta}_{\hat{\beta}^{[k-1]}}) \\ &= \arg \min_{\beta} J_n(\beta, \hat{\eta}_{\hat{\beta}}) + o_P(1) \\ &= \arg \min_{\beta} J_n(\beta, \eta_{\beta_o}) + o_P(1) \\ &= \beta_o + o_P(1). \end{aligned}$$

3.2. Step 2

In Step 2, obtaining $\hat{\eta}_\beta(t)$ for any given values of β and t is equivalent to minimizing $S_n(a; t, \beta)$ in (16). Notice that the dimension $(p + 1)$ of a is typically low, with degrees $p = 0$ or $p = 1$ being the most commonly used in practice. Hence, the minimizer of $S_n(a; t, \beta)$ can be obtained by directly applying the Newton–Raphson iteration: for $k = 0, 1, \dots$,

$$a^{[k+1]}(t; \beta) = a^{[k]}(t; \beta) - \left\{ \frac{\partial^2 S_n(a; t, \beta)}{\partial a \partial a^T} \Big|_{a=a^{[k]}(t; \beta)} \right\}^{-1} \frac{\partial S_n(a; t, \beta)}{\partial a} \Big|_{a=a^{[k]}(t; \beta)},$$

where $a^{[k]}(t; \beta)$ denotes the estimate in the k -th iteration, and:

$$\begin{aligned} \frac{\partial S_n(a; t, \beta)}{\partial a} &= \frac{1}{n} \sum_{i=1}^n p_1(Y_i; X_i^T \beta + \mathbf{t}_i(t)^T a) \mathbf{t}_i(t) w_1(X_i) K_h(T_i - t), \\ \frac{\partial^2 S_n(a; t, \beta)}{\partial a \partial a^T} &= \frac{1}{n} \sum_{i=1}^n p_2(Y_i; X_i^T \beta + \mathbf{t}_i(t)^T a) \mathbf{t}_i(t) \mathbf{t}_i(t)^T w_1(X_i) K_h(T_i - t). \end{aligned}$$

The iterations terminate until the estimate $\hat{\eta}^{[k+1]}(t) = e_{1,p+1}^T a^{[k+1]}(t; \beta)$ stabilizes.

Our numerical studies of the robust-BD estimation indicate that (i) the kernel regression method can be both faster and stabler than the local-linear method; (ii) to estimate the nonparametric component $\eta^o(\cdot)$, the local-linear method outperforms the kernel method, especially at the edges of points $\{T_i\}_{i=1}^n$; (iii) for the performance of the robust estimation of β_o , which is of major interest, there is a relatively negligible difference between choices of using the kernel and local-linear methods in estimating nonparametric components.

4. Asymptotic Property of the Robust-BD Estimators

This section investigates the asymptotic behavior of robust-BD estimators $\hat{\beta}$ and $\hat{\eta}_\beta$, under regularity conditions. The consistency of $\hat{\beta}$ to β_o and uniform consistency of $\hat{\eta}_\beta$ to η^o are given in Theorem 1; the asymptotic normality of $\hat{\beta}$ is obtained in Theorem 2. For the sake of exposition, the asymptotic results will be derived using local-linear estimation with degree $p = 1$. Analogous results can be obtained for local-polynomial methods with lengthier technical details and are omitted.

We assume that $T \in \mathcal{T}$, and let $\mathcal{T}_0 \subseteq \mathcal{T}$ be a compact set. For any continuous function $v : \mathcal{T} \mapsto \mathbb{R}$, define $\|v\|_\infty = \sup_{t \in \mathcal{T}} |v(t)|$ and $\|v\|_{\mathcal{T}_0, \infty} = \sup_{t \in \mathcal{T}_0} |v(t)|$. For a matrix M , the smallest and largest eigenvalues are denoted by $\lambda_j(M)$, $\lambda_{\min}(M)$ and $\lambda_{\max}(M)$, respectively. Let $\|M\| = \sup_{\|x\|=1} \|Mx\| = \{\lambda_{\max}(M^T M)\}^{1/2}$ be the matrix L_2 norm. Denote by \xrightarrow{P} convergence in probability and $\xrightarrow{\mathcal{D}}$ convergence in distribution.

4.1. Consistency

We first present Lemma 1, which states the uniform consistency of $\hat{\eta}_\beta(\cdot)$ to the surrogate function $\eta_\beta(\cdot)$. Theorem 1 gives the consistency of $\hat{\beta}$ and $\hat{\eta}_\beta$.

Lemma 1 (For the non-parametric surrogate $\eta_\beta(\cdot)$). Let $\mathcal{K} \subseteq \mathbb{R}^d$ and $\mathcal{T}_0 \subseteq \mathcal{T}$ be compact sets. Assume Condition A1 and Condition B in the Appendix. If $n \rightarrow \infty$, $h \rightarrow 0$, $nh \rightarrow \infty$, $\log(1/h)/(nh) \rightarrow 0$, then $\sup_{\beta \in \mathcal{K}} \|\hat{\eta}_\beta - \eta_\beta\|_{\mathcal{T}_0, \infty} \xrightarrow{P} 0$.

Theorem 1 (For β_o and $\eta^o(\cdot)$). Assume conditions in Lemma 1.

- (i) If there exists a compact set \mathcal{K}_1 such that $\lim_{n \rightarrow \infty} P(\hat{\beta} \in \mathcal{K}_1) = 1$ and Condition A2 holds, then $\hat{\beta} \xrightarrow{P} \beta_o$.

(ii) Moreover, if Condition A3 holds, then $\|\hat{\eta}_{\hat{\beta}} - \eta^o\|_{\mathcal{T}_0;\infty} \xrightarrow{P} 0$.

4.2. Asymptotic Normality

The asymptotic normality of $\hat{\beta}$ is provided in Theorem 2.

Theorem 2 (For the parametric part β_o). Assume Conditions A and Condition B in the Appendix. If $n \rightarrow \infty$, $nh^4 \rightarrow 0$ and $\log(1/h)/(nh^2) \rightarrow 0$, then:

$$\sqrt{n}(\hat{\beta} - \beta_o) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{H}_0^{-1} \Omega_0^* \mathbf{H}_0^{-1}),$$

where:

$$\mathbf{H}_0 = E \left[p_2(Y; \mathbf{X}^T \beta_o + \eta^o(T)) \left\{ \mathbf{X} + \frac{\partial \eta_{\beta}(T)}{\partial \beta} \Big|_{\beta=\beta_o} \right\} \left\{ \mathbf{X} + \frac{\partial \eta_{\beta}(T)}{\partial \beta} \Big|_{\beta=\beta_o} \right\}^T w_2(\mathbf{X}) \right], \quad (25)$$

and:

$$\begin{aligned} \Omega_0^* = E & \left(p_1^2(Y; \mathbf{X}^T \beta_o + \eta^o(T)) \left[\left\{ \mathbf{X} + \frac{\partial \eta_{\beta}(T)}{\partial \beta} \Big|_{\beta=\beta_o} \right\} w_2(\mathbf{X}) - \frac{\gamma(T)}{g_2(T; T, \beta_o)} w_1(\mathbf{X}) \right] \right. \\ & \left. \times \left[\left\{ \mathbf{X} + \frac{\partial \eta_{\beta}(T)}{\partial \beta} \Big|_{\beta=\beta_o} \right\} w_2(\mathbf{X}) - \frac{\gamma(T)}{g_2(T; T, \beta_o)} w_1(\mathbf{X}) \right]^T \right) \end{aligned} \quad (26)$$

with:

$$\begin{aligned} \gamma(t) &= E \left[p_2(Y; \mathbf{X}^T \beta_o + \eta^o(t)) \left\{ \mathbf{X} + \frac{\partial \eta_{\beta}(t)}{\partial \beta} \Big|_{\beta=\beta_o} \right\} w_2(\mathbf{X}) \mid T = t \right], \\ g_2(t; t, \beta) &= E \{ p_2(Y; \mathbf{X}^T \beta + \eta_{\beta}(t)) w_1(\mathbf{X}) \mid T = t \}. \end{aligned}$$

From Condition A1, (13) and (14), we can show that if $w_1(\cdot) \equiv C w_2(\cdot)$ for some constant $C \in (0, \infty)$, then $\gamma(t) = \mathbf{0}$. In that case, $\Omega_0^* = \Omega_0$, where:

$$\Omega_0 = E \left[p_1^2(Y; \mathbf{X}^T \beta_o + \eta^o(T)) \left\{ \mathbf{X} + \frac{\partial \eta_{\beta}(T)}{\partial \beta} \Big|_{\beta=\beta_o} \right\} \left\{ \mathbf{X} + \frac{\partial \eta_{\beta}(T)}{\partial \beta} \Big|_{\beta=\beta_o} \right\}^T w_2^2(\mathbf{X}) \right]. \quad (27)$$

Consider the conventional PLM in (1), estimated using the classical quadratic loss, identity link and $w_1(\cdot) = w_2(\cdot) \equiv 1$. If $\text{var}(\epsilon \mid \mathbf{X}, T) \equiv \sigma^2$, then $\mathbf{H}_0^{-1} \Omega_0 \mathbf{H}_0^{-1} = \sigma^2 [E\{\text{var}(\mathbf{X} \mid T)\}]^{-1}$, and thus, the result of Theorem 2 agrees with that in [18].

Remark 2. Theorem 2 implies the root- n convergence rate of $\hat{\beta}$. This differs from $\hat{\eta}_{\hat{\beta}}(t)$, which converges at some rate incorporating both the sample size n and the bandwidth h , as seen in the proofs of Lemma 1 and Theorem 2.

5. Robust Inference for β_o Based on BD

In many statistical applications, we will check whether or not a subset of explanatory variables used is statistically significant. Specific examples include:

$$\begin{aligned} H_0 : \beta_{j_0} &= 0, & \text{for } j = j_0, \\ H_0 : \beta_{j_0} &= 0, & \text{for } j = j_1, \dots, j_2. \end{aligned}$$

These forms of linear hypotheses for β_o can be more generally formulated as: (4).

5.1. Wald-Type Test W_n

We propose a robust version of the Wald-type test statistic,

$$W_n = n(\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{g}_0)^T (\mathbf{A}\hat{\mathbf{H}}_0^{-1}\hat{\Omega}_0^*\hat{\mathbf{H}}_0^{-1}\mathbf{A}^T)^{-1} (\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{g}_0), \quad (28)$$

based on the robust-BD estimator $\hat{\boldsymbol{\beta}}$ proposed in Section 2.4, where $\hat{\Omega}_0^*$ and $\hat{\mathbf{H}}_0$ are estimates of Ω_0^* and \mathbf{H}_0 satisfying $\hat{\mathbf{H}}_0^{-1}\hat{\Omega}_0^*\hat{\mathbf{H}}_0^{-1} \xrightarrow{P} \mathbf{H}_0^{-1}\Omega_0^*\mathbf{H}_0^{-1}$. For example,

$$\hat{\mathbf{H}}_0 = \frac{1}{n} \sum_{i=1}^n p_2(Y_i; \mathbf{X}_i^T \hat{\boldsymbol{\beta}} + \hat{\eta}_{\hat{\boldsymbol{\beta}}}(T_i)) \left\{ \mathbf{X}_i + \frac{\partial \hat{\eta}_{\hat{\boldsymbol{\beta}}}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} \right\} \left\{ \mathbf{X}_i + \frac{\partial \hat{\eta}_{\hat{\boldsymbol{\beta}}}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} \right\}^T w_2(\mathbf{X}_i),$$

and:

$$\begin{aligned} \hat{\Omega}_0^* &= \frac{1}{n} \sum_{i=1}^n p_1^2(Y_i; \mathbf{X}_i^T \hat{\boldsymbol{\beta}} + \hat{\eta}_{\hat{\boldsymbol{\beta}}}(T_i)) \left[\left\{ \mathbf{X}_i + \frac{\partial \hat{\eta}_{\hat{\boldsymbol{\beta}}}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} \right\} w_2(\mathbf{X}_i) - \frac{\hat{\gamma}(T_i)}{\hat{g}_2(T_i; T_i, \hat{\boldsymbol{\beta}})} w_1(\mathbf{X}_i) \right] \\ &\quad \times \left[\left\{ \mathbf{X}_i + \frac{\partial \hat{\eta}_{\hat{\boldsymbol{\beta}}}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} \right\} w_2(\mathbf{X}_i) - \frac{\hat{\gamma}(T_i)}{\hat{g}_2(T_i; T_i, \hat{\boldsymbol{\beta}})} w_1(\mathbf{X}_i) \right]^T, \end{aligned}$$

fulfill the requirement, where:

$$\begin{aligned} \frac{\partial \hat{\eta}_{\hat{\boldsymbol{\beta}}}(t)}{\partial \boldsymbol{\beta}} &= - \frac{\sum_{k=1}^n p_2(Y_k; \mathbf{X}_k^T \hat{\boldsymbol{\beta}} + \hat{\eta}_{\hat{\boldsymbol{\beta}}}(t)) \mathbf{X}_k w_1(\mathbf{X}_k) K_h(T_k - t)}{\sum_{k=1}^n p_2(Y_k; \mathbf{X}_k^T \hat{\boldsymbol{\beta}} + \hat{\eta}_{\hat{\boldsymbol{\beta}}}(t)) w_1(\mathbf{X}_k) K_h(T_k - t)}, \\ \hat{\gamma}(t) &= \frac{1}{n} \sum_{k=1}^n p_2(Y_k; \mathbf{X}_k^T \hat{\boldsymbol{\beta}} + \hat{\eta}_{\hat{\boldsymbol{\beta}}}(t)) \left\{ \mathbf{X}_k + \frac{\partial \hat{\eta}_{\hat{\boldsymbol{\beta}}}(t)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} \right\} w_2(\mathbf{X}_k) K_h(T_k - t), \\ \hat{g}_2(t; t, \hat{\boldsymbol{\beta}}) &= \frac{1}{n} \sum_{k=1}^n p_2(Y_k; \mathbf{X}_k^T \hat{\boldsymbol{\beta}} + \hat{\eta}_{\hat{\boldsymbol{\beta}}}(t)) w_1(\mathbf{X}_k) K_h(T_k - t). \end{aligned}$$

Again, we can verify that if $w_1(\cdot) \equiv C w_2(\cdot)$ for some constant $C \in (0, \infty)$ and $\hat{\eta}_{\hat{\boldsymbol{\beta}}}(t)$ is obtained from kernel estimation method, then $\hat{\gamma}(t) = \mathbf{0}$, and hence, $\hat{\Omega}_0^* = \hat{\Omega}_0$, where:

$$\hat{\Omega}_0 = \frac{1}{n} \sum_{i=1}^n p_1^2(Y_i; \mathbf{X}_i^T \hat{\boldsymbol{\beta}} + \hat{\eta}_{\hat{\boldsymbol{\beta}}}(T_i)) \left\{ \mathbf{X}_i + \frac{\partial \hat{\eta}_{\hat{\boldsymbol{\beta}}}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} \right\} \left\{ \mathbf{X}_i + \frac{\partial \hat{\eta}_{\hat{\boldsymbol{\beta}}}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} \right\}^T w_2^2(\mathbf{X}_i).$$

Theorem 3 justifies that under the null, W_n would for large n be distributed as χ_k^2 , thus asymptotically distribution-free.

Theorem 3 (Wald-type test based on robust-BD under H_0). Assume conditions in Theorem 2, and $\hat{\mathbf{H}}_0^{-1}\hat{\Omega}_0^*\hat{\mathbf{H}}_0^{-1} \xrightarrow{P} \mathbf{H}_0^{-1}\Omega_0^*\mathbf{H}_0^{-1}$ in (28). Then, under H_0 in (4), we have that:

$$W_n \xrightarrow{\mathcal{D}} \chi_k^2.$$

Theorem 4 indicates that W_n has a non-trivial local power detecting contiguous alternatives approaching the null at the rate $n^{-1/2}$:

$$H_{1n} : \mathbf{A}\boldsymbol{\beta}_0 - \mathbf{g}_0 = \mathbf{c}/\sqrt{n} \{1 + o(1)\}, \quad (29)$$

where $\mathbf{c} = (c_1, \dots, c_k)^T \neq \mathbf{0}$.

Theorem 4 (Wald-type test based on robust-BD under H_{1n}). Assume conditions in Theorem 2, and $\hat{\mathbf{H}}_0^{-1}\hat{\Omega}_0^*\hat{\mathbf{H}}_0^{-1} \xrightarrow{P} \mathbf{H}_0^{-1}\Omega_0^*\mathbf{H}_0^{-1}$ in (28). Then, under H_{1n} in (29), $W_n \xrightarrow{\mathcal{D}} \chi_k^2(\tau^2)$, where $\tau^2 = \mathbf{c}^T (\mathbf{A}\mathbf{H}_0^{-1}\Omega_0^*\mathbf{H}_0^{-1}\mathbf{A}^T)^{-1} \mathbf{c} > 0$.

To appreciate the discriminating power of W_n in assessing the significance, the asymptotic power is analyzed. Theorem 5 manifests that under the fixed alternative H_1 , $W_n \xrightarrow{P} +\infty$ at the rate n . Thus, W_n has the power approaching to one against fixed alternatives.

Theorem 5 (Wald-type test based on robust-BD under H_1). Assume conditions in Theorem 2, and $\hat{\mathbf{H}}_0^{-1}\hat{\Omega}_0^*\hat{\mathbf{H}}_0^{-1} \xrightarrow{P} \mathbf{H}_0^{-1}\Omega_0^*\mathbf{H}_0^{-1}$ in (28). Then, under H_1 in (4), $n^{-1}W_n \geq \lambda_{\max}^{-1}(\mathbf{A}\mathbf{H}_0^{-1}\Omega_0^*\mathbf{H}_0^{-1}\mathbf{A}^T)\|\mathbf{A}\beta_0 - \mathbf{g}_0\|^2 + o_p(1)$.

For the conventional PLM in (1) estimated using the non-robust quadratic loss, [8] showed the asymptotic equivalence between the Wald-type test and likelihood ratio-type test. Our results in the next Section 5.2 reveal that such equivalence is violated when estimators are obtained using the robust loss functions.

5.2. Likelihood Ratio-Type Test Λ_n

This section explores the degree to which the likelihood ratio-type test is extended to the “robust-BD” for testing the null hypothesis in (4) for the GPLM. The robust-BD test statistic is:

$$\Lambda_n = 2n \left\{ \min_{\beta \in \mathbb{R}^d: \mathbf{A}\beta = \mathbf{g}_0} J_n(\beta, \hat{\eta}_\beta) - J_n(\hat{\beta}, \hat{\eta}_{\hat{\beta}}) \right\}, \quad (30)$$

where $\hat{\beta}$ is the robust-BD estimator for β_0 developed in Section 2.4.

Theorem 6 indicates that the limit distribution of Λ_n under H_0 is a linear combination of independent chi-squared variables, with weights relying on some unknown quantities, thus not distribution free.

Theorem 6 (Likelihood ratio-type test based on robust-BD under H_0). Assume conditions in Theorem 2.

(i) Under H_0 in (4), we obtain:

$$\Lambda_n \xrightarrow{\mathcal{D}} \sum_{j=1}^k \lambda_j \{ (\mathbf{A}\mathbf{H}_0^{-1}\mathbf{A}^T)^{-1} (\mathbf{A}\mathbf{V}_0\mathbf{A}^T) \} Z_j^2,$$

where $\mathbf{V}_0 = \mathbf{H}_0^{-1}\Omega_0^*\mathbf{H}_0^{-1}$ and $\{Z_j\}_{j=1}^k \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$.

(ii) Moreover, if $\psi(r) = r$, $w_1(\mathbf{x}) = w_2(\mathbf{x}) \equiv 1$, and the generating q -function of BD satisfies:

$$q''(m(\mathbf{x}, t)) = -\frac{C}{V(m(\mathbf{x}, t))}, \quad \text{for a constant } C > 0, \quad (31)$$

then under H_0 in (4), we have that $\Lambda_n/C \xrightarrow{\mathcal{D}} \chi_k^2$.

Theorem 7 states that Λ_n has non-trivial local power for identifying contiguous alternatives approaching the null at rate $n^{-1/2}$ and that $\Lambda_n \xrightarrow{P} +\infty$ at the rate n under H_1 , thus having the power approaching to one against fixed alternatives.

Theorem 7 (Likelihood ratio-type test based on robust-BD under H_{1n} and H_1). Assume conditions in Theorem 2. Let $\mathbf{V}_0 = \mathbf{H}_0^{-1}\Omega_0^*\mathbf{H}_0^{-1}$ and $\lambda_j = \lambda_j \{ (\mathbf{A}\mathbf{H}_0^{-1}\mathbf{A}^T)^{-1} (\mathbf{A}\mathbf{V}_0\mathbf{A}^T) \}$, $j = 1, \dots, k$.

- (i) Under H_{1n} in (29), $\Lambda_n \xrightarrow{\mathcal{D}} \sum_{j=1}^k (\sqrt{\lambda_j} Z_j + \mathbf{e}_{j,k}^T \mathbf{S} \mathbf{c})^2$, where $\{Z_j\}_{j=1}^k \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$, and \mathbf{S} is a matrix satisfying $\mathbf{S}^T \mathbf{S} = (\mathbf{A}\mathbf{H}_0^{-1}\mathbf{A}^T)^{-1}$ and $\mathbf{S}(\mathbf{A}\mathbf{V}_0\mathbf{A}^T)\mathbf{S}^T = \text{diag}(\lambda_1, \dots, \lambda_k)$.
- (ii) Under H_1 in (4), $n^{-1}\Lambda_n \geq c\|\mathbf{A}\beta_0 - \mathbf{g}_0\|^2 + o_p(1)$ for a constant $c > 0$.

5.3. Comparison between W_n and Λ_n

In summary, the test W_n has some advantages over the test Λ_n . First, the asymptotic null distribution of W_n is distribution-free, whereas the asymptotic null distribution of Λ_n in general depends on unknown quantities. Second, W_n is invariant with re-scaling the generating q -function of the BD, but Λ_n is not. Third, the computational expense of W_n is much more reduced than that of Λ_n , partly because the integration operations for ρ_q are involved in Λ_n , but not in W_n , and partly because Λ_n requires both unrestricted and restricted parameter estimates, while W_n is useful in cases where restricted parameter estimates are difficult to compute. Thus, W_n will be focused on in numerical studies of Section 6.

6. Simulation Study

We conduct simulation evaluations of the performance of robust-BD estimation methods for general partially linear models. We use the Huber ψ -function $\psi(\cdot)$ with $c = 1.345$. The weight functions are chosen to be $w_1(x) = w_2(x) = 1/\{1 + \sum_{j=1}^d (\frac{x_j - m_j}{s_j})^2\}^{1/2}$, where $x = (x_1, \dots, x_d)^T$, m_j and s_j denote the sample median and sample median absolute deviation of $\{X_{i,j} : i = 1, \dots, n\}$ respectively, $j = 1, \dots, d$. As a comparison, the classical non-robust estimation counterparts correspond to using $\psi(r) = r$ and $w_1(x) = w_2(x) \equiv 1$. Throughout the numerical work, the Epanechnikov kernel function $K(t) = 0.75 \max(1 - t^2, 0)$ is used. All these choices (among many others) are for feasibility; the issues on the trade-off between robustness and efficiency are not pursued further in the paper.

The following setup is used in the simulation studies. The sample size is $n = 200$, and the number of replications is 500. (Incorporating a nonparametric component in the GPLM desires a larger n when the number of covariates increases for better numerical performance.) Local-linear robust-BD estimation is illustrated with the bandwidth parameter h to be 20% of the interval length of the variable T . Results using other data-driven choices of h are similar and are omitted.

6.1. Bernoulli Responses

We generate observations $\{(X_i, T_i, Y_i)\}_{i=1}^n$ randomly from the model,

$$Y \mid (X, T) \sim \text{Bernoulli}(m(X, T)), \quad X \sim N(\mathbf{0}, \Sigma), \quad T \sim \text{Uniform}(0, 1),$$

where $\Sigma = (\sigma_{jk})$ with $\sigma_{jk} = 0.2^{|j-k|}$, and X is independent of T . The link function is $\text{logit}\{m(x, t)\} = x^T \beta_0 + \eta^0(t)$, where $\beta_0 = (2, 2, 0, 0)^T$ and $\eta^0(t) = 2 \sin\{\pi(1 + 2t)\}$. Both the deviance and exponential loss functions are employed as the BD.

For each generated dataset from the true model, we create a contaminated dataset, where 10 data points $(X_{i,j}, Y_i)$ are contaminated as follows: they are replaced by $(X_{i,j}^*, Y_i^*)$, where $Y_i^* = 1 - Y_i$, $i = 1, \dots, 5$,

$$\begin{aligned} X_{1,2}^* &= 5 \text{sign}(U_1 - 0.5), & X_{2,2}^* &= 5 \text{sign}(U_2 - 0.5), & X_{3,2}^* &= 5 \text{sign}(U_3 - 0.5), \\ X_{4,4}^* &= 5 \text{sign}(U_4 - 0.5), & X_{5,1}^* &= 5 \text{sign}(U_5 - 0.5), & X_{6,2}^* &= 5 \text{sign}(U_6 - 0.5), \\ X_{7,3}^* &= 5 \text{sign}(U_7 - 0.5), & X_{8,4}^* &= 5 \text{sign}(U_8 - 0.5), & X_{9,2}^* &= 5 \text{sign}(U_9 - 0.5), \\ X_{10,3}^* &= 5 \text{sign}(U_{10} - 0.5), \end{aligned}$$

with $\{U_i\} \stackrel{\text{i.i.d.}}{\sim} \text{Uniform}(0, 1)$.

Figures 1 and 2 compare the boxplots of $(\hat{\beta}_j - \beta_{j,0})$, $j = 1, \dots, d$, based on the non-robust and robust-BD estimates, where the deviance loss and exponential loss are used as the BD in the top and bottom panels respectively. As seen from Figure 1 in the absence of contamination, both non-robust and robust methods perform comparably well. Besides, the bias in non-robust methods using the exponential loss (with $p_2(y; \theta)$ unbounded) is larger than that of the deviance loss (with $p_2(y; \theta)$ bounded). In the presence of contamination, Figure 2 reveals that the robust method is more effective in decreasing the estimation bias without excessively increasing the estimation variance.

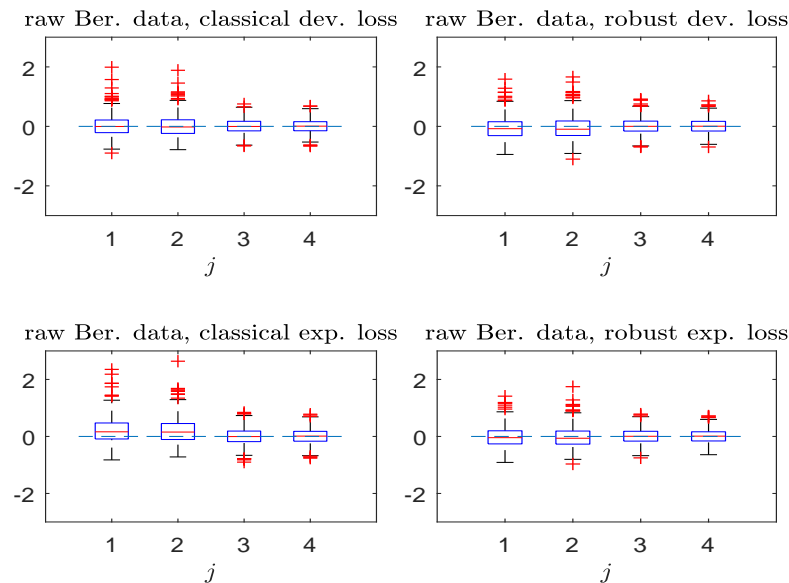


Figure 1. Simulated Bernoulli response data without contamination. Boxplots of $(\hat{\beta}_j - \beta_{j0})$, $j = 1, \dots, d$ (from left to right). (**Left panels**): non-robust method; (**right panels**): robust method.

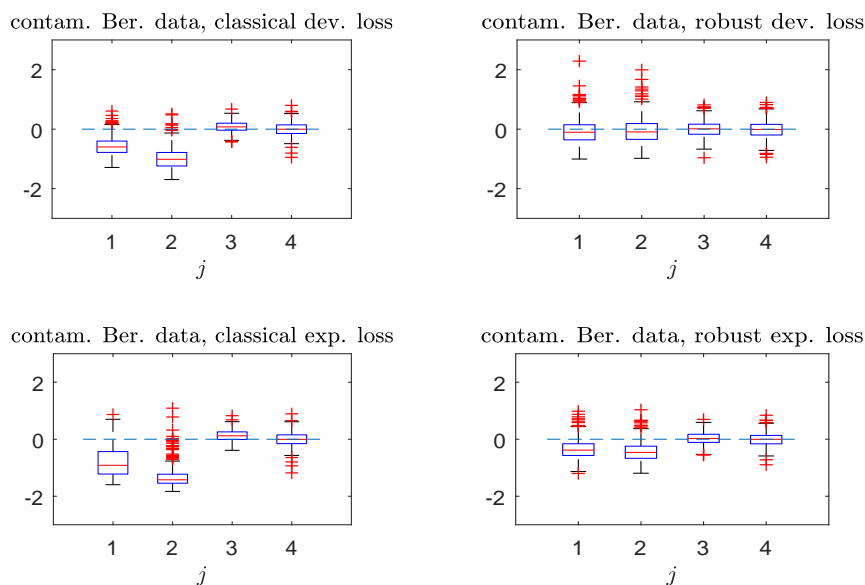


Figure 2. Simulated Bernoulli response data with contamination. The captions are identical to those in Figure 1.

For each replication, we calculate $\text{MSE}(\hat{\eta}) = n^{-1} \sum_{i=1}^n \{\hat{\eta}_{\hat{\beta}}(t_i) - \eta^0(t_i)\}^2$. Figures 3 and 4 compare the plots of $\hat{\eta}_{\hat{\beta}}(t)$ from typical samples, using non-robust and robust-BD estimates, where the deviance loss and exponential loss are used as the BD in the top and bottom panels, respectively. There, the typical sample in each panel is selected in a way such that its MSE value corresponds to the 50-th percentile among the MSE-ranked values from 500 replications. These fitted curves reveal little difference between using the robust and non-robust methods, in the absence of contamination. For contaminated cases, robust estimates perform slightly better than non-robust estimates. Moreover, the boundary bias issue arising from the curve estimates at the edges using the local constant method can be ameliorated by using the local-linear method.

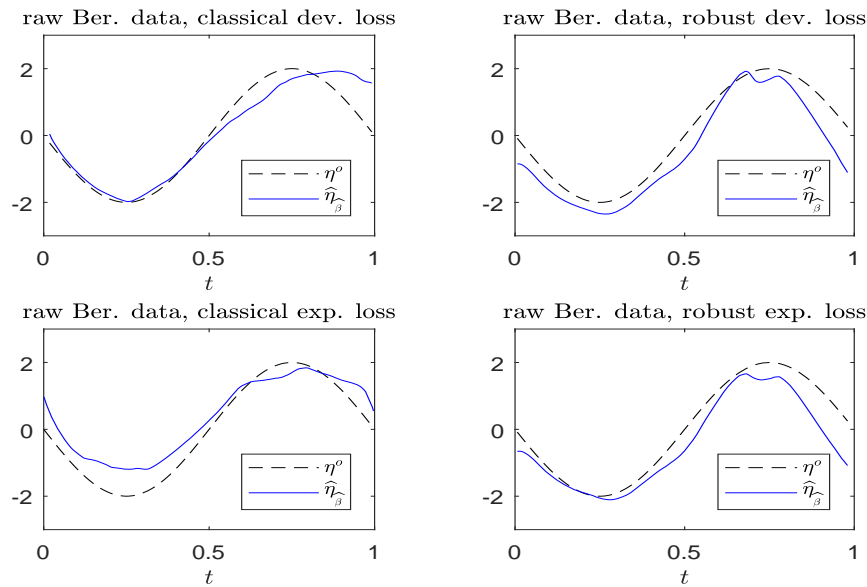


Figure 3. Simulated Bernoulli response data without contamination. Plots of $\eta^0(t)$ and $\hat{\eta}_{\hat{\beta}}(t)$. (Left panels): non-robust method; (right panels): robust method.

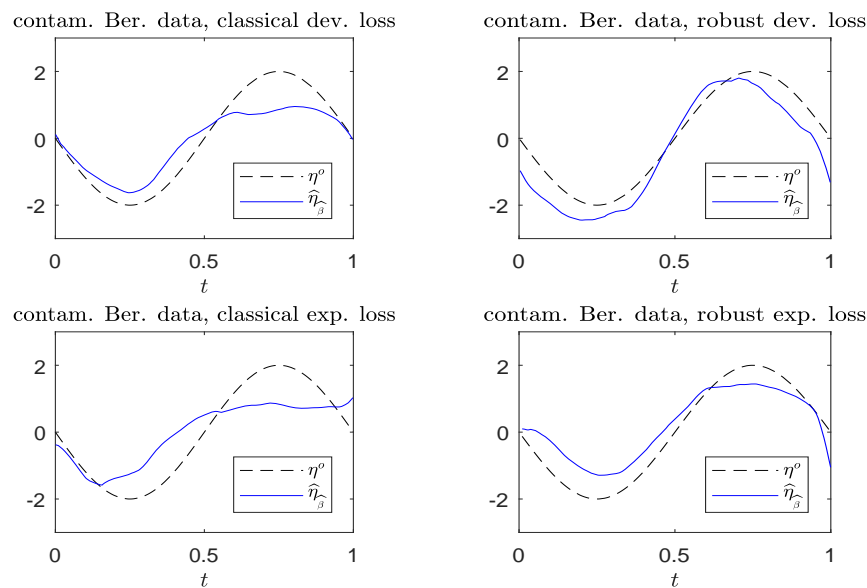


Figure 4. Simulated Bernoulli response data with contamination. Plots of $\eta^0(t)$ and $\hat{\eta}_{\hat{\beta}}(t)$. (Left panels): non-robust method; (right panels): robust method.

6.2. Gaussian Responses

We generate independent observations $\{(\mathbf{X}_i, T_i, Y_i)\}_{i=1}^n$ from (\mathbf{X}, T, Y) satisfying:

$$Y \mid (\mathbf{X}, T) \sim N(m(\mathbf{X}, T), \sigma^2), \quad (\mathbf{X}, \Phi^{-1}(T)) \sim N(\mathbf{0}, \Sigma),$$

where $\sigma = 1$, $\Sigma = (\sigma_{jk})$ with $\sigma_{jk} = 0.2^{|j-k|}$, Φ denotes the CDF of the standard normal distribution. The link function is $m(\mathbf{x}, t) = \mathbf{x}^T \boldsymbol{\beta}_0 + \eta^0(t)$, where $\boldsymbol{\beta}_0 = (2, -2, 1, -1, 0, 0)^T$ and $\eta^0(t) = 2 \sin\{\pi(1 + 2t)\}$. The quadratic loss is utilized as the BD.

For each dataset simulated from the true model, a contaminated data-set is created, where 10 data points $(X_{i,j}, Y_i)$ are subject to contamination. They are replaced by $(X_{i,j}^*, Y_i^*)$, where $Y_i^* = Y_i I\{|Y_i - m(X_i, T_i)|/\sigma > 2\} + 15 I\{|Y_i - m(X_i, T_i)|/\sigma \leq 2\}$, $i = 1, \dots, 10$,

$$\begin{aligned} X_{1,2}^* &= 5 \operatorname{sign}(U_1 - 0.5), & X_{2,2}^* &= 5 \operatorname{sign}(U_2 - 0.5), & X_{3,2}^* &= 5 \operatorname{sign}(U_3 - 0.5), \\ X_{4,4}^* &= 5 \operatorname{sign}(U_4 - 0.5), & X_{5,6}^* &= 5 \operatorname{sign}(U_5 - 0.5), & X_{6,1}^* &= 5 \operatorname{sign}(U_6 - 0.5), \\ X_{7,2}^* &= 5 \operatorname{sign}(U_7 - 0.5), & X_{8,3}^* &= 5 \operatorname{sign}(U_8 - 0.5), & X_{9,4}^* &= 5 \operatorname{sign}(U_9 - 0.5), \\ X_{10,5}^* &= 5 \operatorname{sign}(U_{10} - 0.5), \end{aligned}$$

with $\{U_i\} \stackrel{\text{i.i.d.}}{\sim} \text{Uniform}(0, 1)$.

Figures 5 and 6 compare the boxplots of $(\hat{\beta}_j - \beta_{j,0})$, $j = 1, \dots, d$, on the top panels, and plots of $\hat{\eta}_{\hat{\beta}}(t)$ from typical samples, on the bottom panels, using the non-robust and robust-BD estimates. The typical samples are selected similar to those in Section 6.1. The simulation results in Figure 5 indicate that the robust method performs, as well as the non-robust method for estimating both the parameter vector and non-parametric curve in non-contaminated cases. Figure 6 reveals that the robust estimates are less sensitive to outliers than the non-robust counterparts. Indeed, the non-robust method yields a conceivable bias for parametric estimation, and non-parametric estimation is worse than that of the robust method.

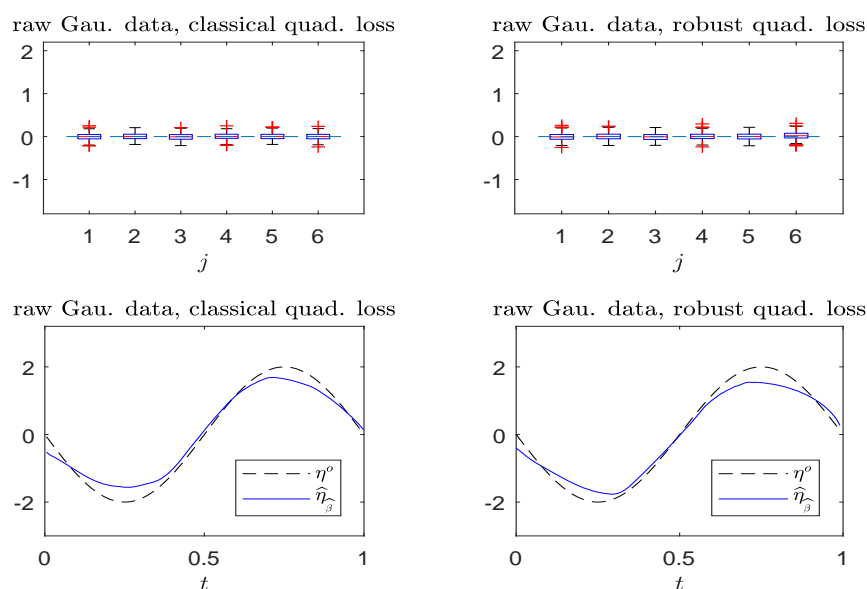


Figure 5. Simulated Gaussian response data without contamination. Top panels: boxplots of $(\hat{\beta}_j - \beta_{j,0})$, $j = 1, \dots, d$ (from left to right). Bottom panels: plots of $\eta^o(t)$ and $\hat{\eta}_{\hat{\beta}}(t)$. (**Left panels**): non-robust method; (**right panels**): robust method.

Figure 7 gives the QQ plots of the (first to 95-th) percentiles of the Wald-type statistic W_n versus those of the χ_2^2 distribution for testing the null hypothesis:

$$H_0 : \beta_{2,0} = -2 \text{ and } \beta_{4,0} = -1. \quad (32)$$

The plots depict that in both clean and contaminated cases, the robust W_n (in right panels) closely follows the χ_2^2 distribution, lending support to Theorem 3. On the other hand, the non-robust W_n agrees well with the χ_2^2 distribution in clean data; the presence of a small number of outlying data points severely distorts the sampling distribution of the non-robust W_n (in the bottom left panel) from the χ_2^2 distribution, yielding inaccurate levels of the test.

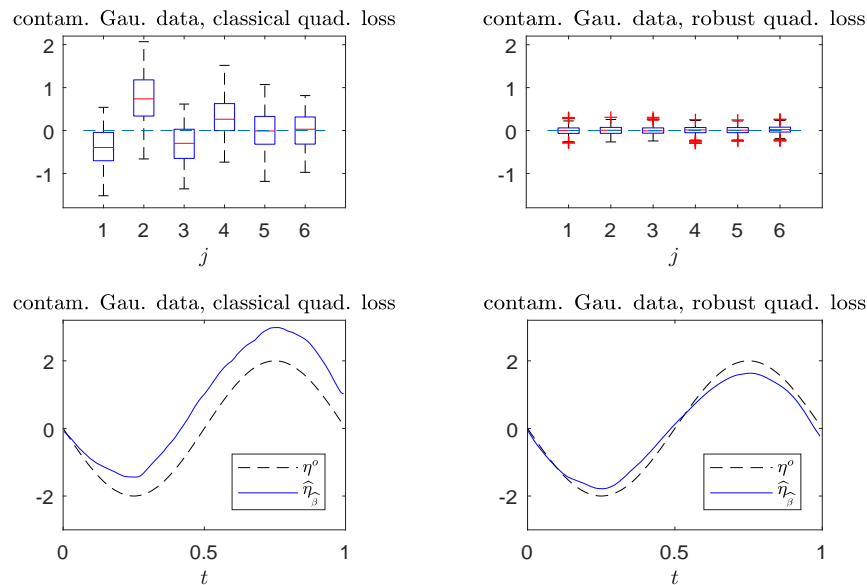


Figure 6. Simulated Gaussian response data with contamination. Top panels: boxplots of $(\hat{\beta}_j - \beta_{j0})$, $j = 1, \dots, d$ (from left to right). Bottom panels: plots of $\eta^o(t)$ and $\hat{\eta}_{\hat{\beta}}(t)$. (**Left panels**): non-robust method; (**right panels**): robust method.

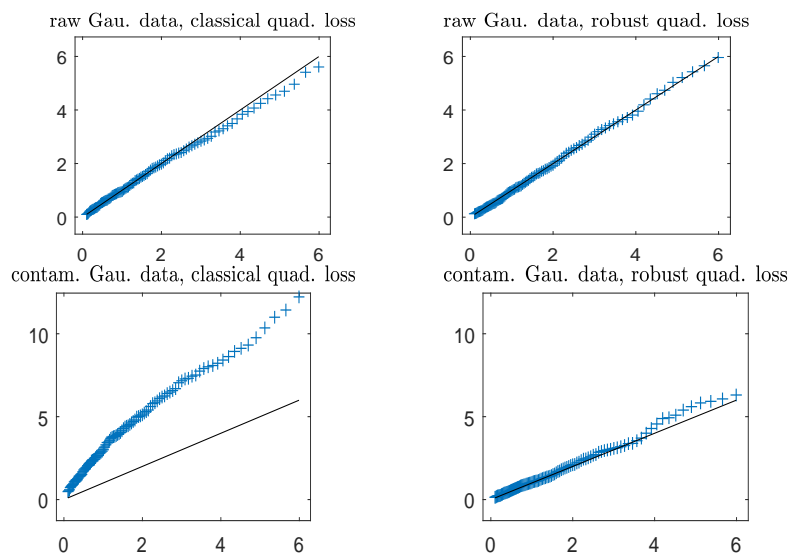


Figure 7. Simulated Gaussian response data with contamination. Empirical quantiles (on the y -axis) of the Wald-type statistics W_n versus quantiles (on the x -axis) of the χ^2 distribution. Solid line: the 45 degree reference line. (**Left panels**): non-robust method; (**right panels**): robust method.

To assess the stability of the power of the Wald-type test for testing the hypothesis (32), we evaluate the power in a sequence of alternatives with parameters $\beta_0 + \Delta c$ for each given Δ , where $c = \beta_0 + (1, \dots, 1)^T$. Figure 8 plots the empirical rejection rates of the null model in the non-contaminated case and the contaminated case. The price to pay for the robust W_n is a little loss of power in the non-contaminated cases. However, under contamination, a very different behavior is observed. The observed power curve of the robust W_n is close to those attained in the non-contaminated case. On the contrary, the non-robust W_n is less informative, since its power curve is much lower than that of the robust W_n against the alternative hypotheses with $\Delta \neq 0$, but higher than the nominal level at the null hypothesis with $\Delta = 0$.

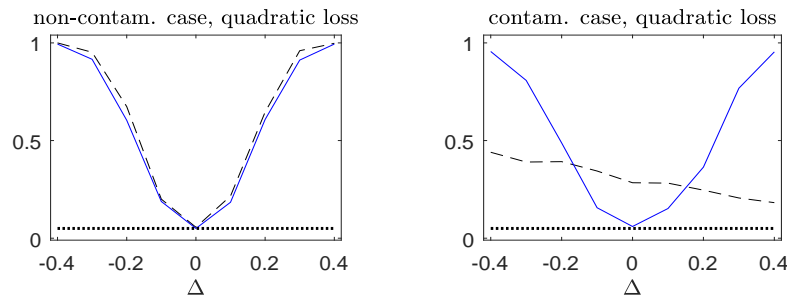


Figure 8. Observed power curves of tests for the Gaussian response data. The dashed line corresponds to the non-robust Wald-type test W_n ; the solid line corresponds to the robust W_n ; the dotted line indicates the 5% nominal level. (**Left panels**): non-contaminated case; (**right panels**): contaminated case.

7. Real Data Analysis

Two real datasets are analyzed. In both cases, the quadratic loss is set to be the BD, and the nonparametric function is fitted via local-linear regression method, where the bandwidth parameter is chosen to be 25% of the interval length of the variable T . Choices of the Huber ψ -function and weight functions are identical to those in Section 6.

7.1. Example 1

The dataset studied in [19] consists of 2447 observations on three variables, $\log(\text{wage})$, age and education, for women. It is of interest to learn how wages change with years of age and years of education. It is anticipated to find an increasing regression function of $Y = \log(\text{wage})$ in $T = \text{age}$ as well as in $X_1 = \text{education}$. We fit a partially linear model $Y = \eta(T) + \beta_1 X_1 + \epsilon$. Profiles of the fitted nonparametric functions $\hat{\eta}(\cdot)$ in Figure 9 indeed exhibit the overall upward trend in age. The coefficient estimate is $\hat{\beta}_1 = 0.0809$ with standard error 0.0042 using the non-robust method, and is $\hat{\beta}_1 = 0.1334$ with standard error 0.0046 by means of the robust method. It is seen that robust estimates are similar to the non-robust counterparts. Our evaluation, based on both the non-robust and robust methods, supports the predicted result in theoretical and empirical literature in socio-economical studies.

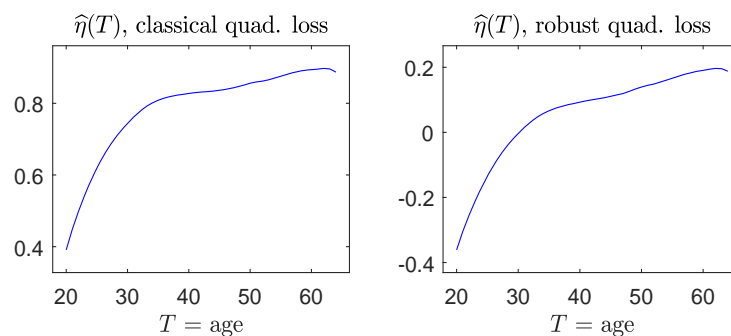


Figure 9. The dataset in [19]. (**Left panels**): estimate of $\eta(T)$ via the non-robust quadratic loss; (**right panels**): estimate of $\eta(T)$ via the robust quadratic loss.

7.2. Example 2

We analyze an employee dataset (Example 11.3 of [20]) of the Fifth National Bank of Springfield, based on year 1995 data. The bank, whose name has been changed, was charged in court with that its female employees received substantially smaller salaries than its male employees. For each of its 208 employees, the dataset consists of seven variables, EducLev (education level), JobGrade (job grade), YrHired (year that an employee was hired), YrBorn (year that an employee was born), Female

(indicator of being female), YrsPrior (years of work experience at another bank before working at the Fifth National bank), and Salary (current annual salary in thousands of dollars).

To explain variation in salary, we fit a partial linear model, $Y = \eta(T) + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \epsilon$, for $Y = \log(\text{Salary})$, $T = \text{Age}$, $X_1 = \text{Female}$, $X_2 = \text{YrHired}$, $X_3 = \text{EducLev}$, $X_4 = \text{JobGrade}$ and $X_5 = \text{YrsPrior}$, where $\text{Age} = 95 - \text{YrBorn}$ is age. Table 1 presents parameter estimates and their standard errors (given within brackets), along with p -values calculated from the Wald-type test W_n . Figure 10 depicts the estimated nonparametric functions.

It is interesting to note that for this dataset, results from using the robust and non-robust methods make a difference in drawing conclusions. For example, from Table 1, the non-robust method gives the estimate of parameter β_1 for gender to be below zero, which may be interpreted as the evidence of discrimination against female employees in salary and lends support to the plaintiff. In contrast, the robust method yields $\hat{\beta}_1 > 0$, which does not indicate that gender has an adverse effect. (A similar conclusion made from penalized-likelihood was obtained in Section 4.1 of [21]). Moreover, the estimated nonparametric functions $\hat{\eta}(\cdot)$ obtained from non-robust and robust methods are qualitatively different: the former method does not deliver a monotone increasing pattern with Age, whereas the latter method does. Whether or not the difference was caused by outlying observations will be an interesting issue to be investigated.

Table 1. Parameter estimates and p -values for partially linear model of the dataset in [20]

Variable	Classical-BD Estimation			Robust-BD Estimation		
	Estimate (s.e.)	p -Value of W_n		Estimate (s.e.)	p -Value of W_n	
Female	−0.0491 (0.0232)	0.0339		0.0530 (0.0323)	0.1010	
YrHired	−0.0093 (0.0026)	0.0005		0.0359 (0.0086)	0.0000	
EducLev	0.0179 (0.0079)	0.0228		−0.0133 (0.0131)	0.3103	
JobGrade	0.0899 (0.0075)	0.0000		0.1672 (0.0168)	0.0000	
YrsPrior	0.0033 (0.0023)	0.1528		−0.0050 (0.0061)	0.4104	

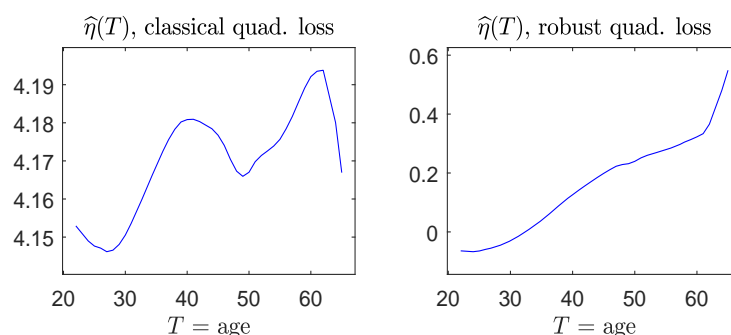


Figure 10. The dataset in [20]. (Left panel): estimate of $\eta(T)$ via the non-robust quadratic loss; (right panel): estimate of $\eta(T)$ via the robust quadratic loss.

8. Discussion

Over the past two decades, nonparametric inference procedures for testing hypotheses concerning nonparametric regression functions have been developed extensively. See [22–26] and the references therein. The work on the generalized likelihood ratio test [24] offers light into nonparametric inference, based on function estimation under nonparametric models, using the quadratic loss function as the error measure. These works do not directly deal with the robust procedure. Exploring the inference on nonparametric functions, such as $\eta^o(t)$ in GPLM associated with a scalar variable T and the additive structure $\sum_{d=1}^D \eta_d^o(t_d)$ as in [27] with a vector variable $T = (T_1, \dots, T_D)$, estimated via the “robust-BD” as the error measure, when there are possible outlying data points, will be the future work.

This paper utilizes the class BD of loss functions, the optimal choice of which depends on specific settings and criteria. For e.g., regression and classification will utilize different loss functions, and thus further study on optimality is desirable.

Some recent work on partially linear models in econometrics includes [28–30]. There, the nonparametric function is approximated via linear expansions, with the number of coefficients diverging with n . Developing inference procedures to be resistant to outliers could be of interest.

Acknowledgments: The authors thank the two referees for insightful comments and suggestions. The research is supported by the U.S. NSF Grants DMS–1712418, DMS–1505367, CMMI–1536978, DMS–1308872, the Wisconsin Alumni Research Foundation and the National Natural Science Foundation of China, grants 11690014.

Author Contributions: C.Z. conceived and designed the experiments; C.Z. analyzed the data; Z.Z. contributed to discussions and analysis tools; C.Z. wrote the paper.

Conflicts of Interest: The authors declare no conflict of interest.

Appendix A. Proofs of Main Results

Throughout the proof, C represents a generic finite constant. We impose some regularity conditions, which may not be the weakest, but facilitate the technical derivations.

Notation:

For integers $j \geq 0$, $\mu_j(K) = \int u^j K(u) du$; $\mathbf{c}_p = (\mu_{p+1}(K), \dots, \mu_{2p+1}(K))^T$; $\mathcal{S} = (\mu_{j+k-2}(K))_{1 \leq j, k \leq p+1}$. Define: $\eta(\mathbf{x}, t) = F(m(\mathbf{x}, t)) = \mathbf{x}^T \boldsymbol{\beta}_0 + \eta^0(t)$; $\eta_i = \eta(\mathbf{X}_i, T_i)$. Set $\eta_i(t; \boldsymbol{\beta}) = \mathbf{X}_i^T \boldsymbol{\beta} + \eta_{\beta}(t) + \sum_{k=1}^p (T_i - t)^k \eta_{\beta}^{(k)}(t) / k!$; $g_1(\tau; t, \boldsymbol{\beta}) = E\{p_1(Y_i; \eta_i(t; \boldsymbol{\beta})) w_1(\mathbf{X}_i) \mid T_i = \tau\}$; $g_2(\tau; t, \boldsymbol{\beta}) = E\{p_2(Y_i; \eta_i(t; \boldsymbol{\beta})) w_1(\mathbf{X}_i) \mid T_i = \tau\}$.

Condition A:

- A1. $\eta_{\beta}(t)$ is the unique minimizer of $S(a; t, \boldsymbol{\beta})$ with respect to $a \in \mathbb{R}^1$.
- A2. $\boldsymbol{\beta}_0 \in \mathbb{R}^d$ is the unique minimizer of $J(\boldsymbol{\beta}, \eta_{\beta})$ with respect to $\boldsymbol{\beta}$, where $d \geq 1$.
- A3. $\eta^0(\cdot) = \eta_{\beta_0}(\cdot)$.

Condition B:

- B1. The function $\rho_q(y, \mu)$ is continuous and bounded. The functions $p_1(y; \theta)$, $p_2(y; \theta)$, $p_3(y; \theta)$, $w_1(\cdot)$ and $w_2(\cdot)$ are bounded; $p_2(y; \theta)$ is continuous in θ .
- B2. The kernel function K is Lipschitz continuous, a symmetric probability density function with bounded support. The matrix \mathcal{S} is positive definite.
- B3. The marginal density $f_T(t)$ of T is a continuous function, uniformly bounded away from zero and ∞ for $t \in \mathcal{T}_0$.
- B4. The function $S(a; t, \boldsymbol{\beta})$ is continuous and $\eta_{\beta}(t)$ is a continuous function of $(t, \boldsymbol{\beta})$.
- B5. Assume $g_2(\tau; t, \boldsymbol{\beta})$ is continuous in τ ; $g_2(t; t, \boldsymbol{\beta})$ is continuous in $t \in \mathcal{T}_0$.
- B6. Functions $\eta_{\beta}(t)$ and $\eta^0(t)$ are $(p+1)$ -times continuously differentiable at t .
- B7. The link function $F(\cdot)$ is monotone increasing and a bijection, $F^{(3)}(\cdot)$ is continuous, and $F^{(1)}(\cdot) > 0$. The matrix $\text{var}(\mathbf{X} \mid T = t)$ is positive definite for a.e. t .
- B8. The matrix \mathbf{H}_0 in (25) is invertible; Ω_0^* in (26) is positive-definite.
- B9. $\hat{\eta}_{\beta}(t)$ and $\eta_{\beta}(t)$ are continuously differentiable with respect to $(t, \boldsymbol{\beta})$, and twice continuously differentiable with respect to $\boldsymbol{\beta}$ such that for any $1 \leq j, k \leq d$, $\frac{\partial^2}{\partial \beta_j \partial \beta_k} \eta_{\beta}(t) |_{\boldsymbol{\beta}=\boldsymbol{\beta}_0}$ is bounded. Furthermore, for any $1 \leq j, k \leq d$, $\frac{\partial^2}{\partial \beta_j \partial \beta_k} \eta_{\beta}(t)$ satisfies the equicontinuity condition:

$$\forall \varepsilon > 0, \exists \delta_{\varepsilon} > 0 : \|\boldsymbol{\beta}_1 - \boldsymbol{\beta}_0\| < \delta_{\varepsilon} \implies \left\| \frac{\partial^2}{\partial \beta_j \partial \beta_k} \eta_{\beta} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_1} - \frac{\partial^2}{\partial \beta_j \partial \beta_k} \eta_{\beta} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\|_{\infty} < \varepsilon.$$

Note that Conditions A, B2–B5 and B8–B9 were similarly used in [9]. Conditions B1 and B7 follow [10]. Condition B6 is due to the local p -th-degree polynomial regression estimation.

Proof of Lemma 1: From Condition A1, we obtain $E\{p_1(Y; X^T \beta + \eta_\beta(t))w_1(X) \mid T = t\} = 0$ and $E\{p_2(Y; X^T \beta + \eta_\beta(t))w_1(X) \mid T = t\} > 0$, i.e.,

$$g_1(t; t, \beta) = E\{p_1(Y; X^T \beta + \eta_\beta(t))w_1(X) \mid T = t\} = 0, \quad (A1)$$

$$g_2(t; t, \beta) = E\{p_2(Y; X^T \beta + \eta_\beta(t))w_1(X) \mid T = t\} > 0. \quad (A2)$$

Define by $\eta_\beta^{(0, \dots, p)}(t) = (\eta_\beta(t), \eta_\beta^{(1)}(t), \dots, \eta_\beta^{(p)}(t)/p!)^T$ the vector of $\eta_\beta(t)$ along with re-scaled derivatives with respect to t up to the order p . Note that:

$$\begin{aligned} \eta_i(t; \beta) &= X_i^T \beta + \sum_{k=0}^p (T_i - t)^k \frac{\eta_\beta^{(k)}(t)}{k!} \\ &= X_i^T \beta + \mathbf{t}_i(t)^T \boldsymbol{\eta}_\beta^{(0, \dots, p)}(t) \\ &= X_i^T \beta + \{H^{-1} \mathbf{t}_i(t)\}^T H \boldsymbol{\eta}_\beta^{(0, \dots, p)}(t) \\ &= X_i^T \beta + \mathbf{t}_i^*(t)^T H \boldsymbol{\eta}_\beta^{(0, \dots, p)}(t), \end{aligned}$$

where $H = \text{diag}\{(1, h, \dots, h^p)\}$ and $\mathbf{t}_i^*(t) = H^{-1} \mathbf{t}_i(t) = (1, (T_i - t)/h, \dots, (T_i - t)^p/h^p)^T$ denotes the re-scaled $\mathbf{t}_i(t)$. Then:

$$\begin{aligned} &X_i^T \beta + \mathbf{t}_i(t)^T \mathbf{a} \\ &= X_i^T \beta + \mathbf{t}_i^*(t)^T H \mathbf{a} \\ &= X_i^T \beta + \mathbf{t}_i^*(t)^T H \boldsymbol{\eta}_\beta^{(0, \dots, p)}(t) + \mathbf{t}_i^*(t)^T H \{\mathbf{a} - \boldsymbol{\eta}_\beta^{(0, \dots, p)}(t)\} \\ &= \eta_i(t; \beta) + \mathbf{t}_i^*(t)^T H \{\mathbf{a} - \boldsymbol{\eta}_\beta^{(0, \dots, p)}(t)\}. \end{aligned}$$

Hence, we rewrite (16) as:

$$S_n(\mathbf{a}; t, \beta) = \frac{1}{n} \sum_{i=1}^n \rho_q(Y_i, F^{-1}(\eta_i(t; \beta) + \mathbf{t}_i^*(t)^T H \{\mathbf{a} - \boldsymbol{\eta}_\beta^{(0, \dots, p)}(t)\})) w_1(X_i) K_h(T_i - t).$$

Therefore, $\hat{\mathbf{a}}(t, \beta)$ minimizing $S_n(\mathbf{a}; t, \beta)$ is equivalent to the one minimizing:

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n \left\{ \rho_q(Y_i, F^{-1}(\eta_i(t; \beta) + \mathbf{t}_i^*(t)^T H \{\mathbf{a} - \boldsymbol{\eta}_\beta^{(0, \dots, p)}(t)\})) \right. \\ &\quad \left. - \rho_q(Y_i, F^{-1}(\eta_i(t; \beta))) \right\} w_1(X_i) K_h(T_i - t) \end{aligned}$$

with respect to \mathbf{a} . It follows that $\hat{\mathbf{a}}^*(t, \beta)$, defined by $\hat{\mathbf{a}}^*(t, \beta) = \sqrt{nh} H \{\hat{\mathbf{a}}(t, \beta) - \boldsymbol{\eta}_\beta^{(0, \dots, p)}(t)\}$, minimizes:

$$\begin{aligned} G_n(\mathbf{a}^*; t, \beta) &= nh \left[\frac{1}{n} \sum_{i=1}^n \left\{ \rho_q(Y_i, F^{-1}(\eta_i(t; \beta) + \{a_n \mathbf{t}_i^*(t)^T \mathbf{a}^*\})) - \rho_q(Y_i, F^{-1}(\eta_i(t; \beta))) \right\} \right. \\ &\quad \left. w_1(X_i) K_h(T_i - t) \right] \end{aligned}$$

with respect to $\mathbf{a}^* \in \mathbb{R}^{p+1}$, where $a_n = 1/\sqrt{nh}$. Note that for any fixed \mathbf{a}^* , $|\mathbf{t}_i^*(t)^T \mathbf{a}^*| \leq C$. By Taylor expansion,

$$\begin{aligned} G_n(\mathbf{a}^*; t, \beta) &= nh \left(a_n \left[\frac{1}{n} \sum_{i=1}^n p_1(Y_i; \eta_i(t; \beta)) \{\mathbf{t}_i^*(t)^T \mathbf{a}^*\} w_1(X_i) K_h(T_i - t) \right] \right. \\ &\quad \left. + a_n^2 \frac{1}{2} \left[\frac{1}{n} \sum_{i=1}^n p_2(Y_i; \eta_i(t; \beta)) \{\mathbf{t}_i^*(t)^T \mathbf{a}^*\}^2 w_1(X_i) K_h(T_i - t) \right] \right) \end{aligned}$$

$$\begin{aligned}
& + a_n^3 \frac{1}{6} \left[\frac{1}{n} \sum_{i=1}^n p_3(Y_i; \eta_i^*(t; \beta)) \{ \mathbf{t}_i^*(t)^T \mathbf{a}^* \}^3 w_1(\mathbf{X}_i) K_h(T_i - t) \right] \\
& = I_{n,1} + I_{n,2} + I_{n,3},
\end{aligned}$$

where $\eta_i^*(t; \beta)$ is located between $\eta_i(t; \beta)$ and $\eta_i(t; \beta) + \{a_n \mathbf{t}_i^*(t)^T \mathbf{a}^*\}$. We notice that:

$$I_{n,1} \equiv \sqrt{nh} \mathbf{W}_n(t, \beta)^T \mathbf{a}^*,$$

where:

$$\mathbf{W}_n(t, \beta) = \frac{1}{n} \sum_{i=1}^n p_1(Y_i; \eta_i(t; \beta)) \mathbf{t}_i^*(t) w_1(\mathbf{X}_i) K_h(T_i - t);$$

also, Lemma A1 implies:

$$\begin{aligned}
I_{n,2} &= n h a_n^2 \frac{1}{2} \mathbf{a}^{*T} \left[\frac{1}{n} \sum_{i=1}^n p_2(Y_i; \eta_i(t; \beta)) \{ \mathbf{t}_i^*(t) \mathbf{t}_i^*(t)^T \} w_1(\mathbf{X}_i) K_h(T_i - t) \right] \mathbf{a}^* \\
&= \frac{1}{2} \mathbf{a}^{*T} \mathbf{S}_2(t, \beta) \mathbf{a}^* + o_p(1),
\end{aligned}$$

where:

$$\mathbf{S}_2(t, \beta) = g_2(t; t, \beta) f_T(t) \mathcal{S} \succ \mathbf{0}$$

by (A2), Condition B2 and B5; and (by using $X_n = O_p(E(|X_n|))$):

$$I_{n,3} \leq C O_p(n h a_n^3) = O_p(1/\sqrt{nh}) = o_p(1).$$

Then:

$$G_n(\mathbf{a}^*; t, \beta) = \sqrt{nh} \mathbf{W}_n(t, \beta)^T \mathbf{a}^* + \frac{1}{2} \mathbf{a}^{*T} \mathbf{S}_2(t, \beta) \mathbf{a}^* + o_p(1),$$

where $\mathbf{a}^{*T} \mathbf{S}_2(t, \beta) \mathbf{a}^* = (\mathbf{a}^{*T} \mathcal{S} \mathbf{a}^*) g_2(t; t, \beta) f_T(t)$ is continuous in $t \in \mathcal{T}_0$ by B3 and B5.

We now examine $\mathbf{W}_n(t, \beta)$. Note that:

$$\begin{aligned}
\text{var}\{\mathbf{W}_n(t, \beta)\} &= \frac{1}{n} \text{var}\{p_1(Y_i; \eta_i(t; \beta)) \mathbf{t}_i^*(t) w_1(\mathbf{X}_i) K_h(T_i - t)\} \\
&\leq \frac{1}{n} E\left[p_1^2(Y_i; \eta_i(t; \beta)) \{ \mathbf{t}_i^*(t) \mathbf{t}_i^*(t)^T \} w_1^2(\mathbf{X}_i) \{K_h(T_i - t)\}^2\right] \\
&\leq \frac{C}{n} E\left[\frac{1}{h^2} \left\{K\left(\frac{T_i - t}{h}\right)\right\}^2\right] \\
&= \frac{C}{nh}.
\end{aligned}$$

To evaluate $E\{\mathbf{W}_n(t, \beta)\}$, it is easy to see that for each $j \in \{0, 1, \dots, p\}$,

$$\begin{aligned}
e_{j+1, p+1}^T E\{\mathbf{W}_n(t, \beta)\} &= E\{p_1(Y_i; \eta_i(t; \beta)) e_{j+1, p+1}^T \mathbf{t}_i^*(t) w_1(\mathbf{X}_i) K_h(T_i - t)\} \\
&= E\left\{p_1(Y_i; \eta_i(t; \beta)) \left(\frac{T_i - t}{h}\right)^j w_1(\mathbf{X}_i) K_h(T_i - t)\right\} \\
&= E\left[E\{p_1(Y_i; \eta_i(t; \beta)) w_1(\mathbf{X}_i) \mid T_i\} \left(\frac{T_i - t}{h}\right)^j K_h(T_i - t)\right] \\
&= E\left\{g_1(T_i; t, \beta) \left(\frac{T_i - t}{h}\right)^j K_h(T_i - t)\right\} \\
&= \int g_1(y; t, \beta) \left(\frac{y - t}{h}\right)^j \frac{1}{h} K\left(\frac{y - t}{h}\right) f_T(y) dy \\
&= \int g_1(t + hx; t, \beta) x^j K(x) f_T(t + hx) dx.
\end{aligned}$$

Note that by Taylor expansion,

$$\eta_{\beta}(t+hx) = \sum_{k=0}^p (hx)^k \frac{\eta_{\beta}^{(k)}(t)}{k!} + (hx)^{p+1} \frac{\eta_{\beta}^{(p+1)}(t)}{(p+1)!} + o(h^{p+1}).$$

This combined with the facts (A1) and (A2) give that:

$$\begin{aligned} & g_1(t+hx; t, \beta) \\ = & \mathbb{E} \left\{ p_1 \left(Y; \mathbf{X}^T \beta + \sum_{k=0}^p (hx)^k \frac{\eta_{\beta}^{(k)}(t)}{k!} \right) w_1(\mathbf{X}) \mid T = t+hx \right\} \\ = & \mathbb{E} \left[p_1(Y; \mathbf{X}^T \beta + \eta_{\beta}(t+hx)) w_1(\mathbf{X}) \right. \\ & \left. + p_2(Y; \mathbf{X}^T \beta + \eta_{\beta}(t+hx)) \left\{ \sum_{k=0}^p (hx)^k \frac{\eta_{\beta}^{(k)}(t)}{k!} - \eta_{\beta}(t+hx) \right\} w_1(\mathbf{X}) \mid T = t+hx \right] \\ & + o(h^{p+1}) \\ = & g_1(t+hx; t+hx, \beta) - (hx)^{p+1} \frac{\eta_{\beta}^{(p+1)}(t)}{(p+1)!} g_2(t+hx; t+hx, \beta) + o(h^{p+1}) \\ = & -(hx)^{p+1} \frac{\eta_{\beta}^{(p+1)}(t)}{(p+1)!} g_2(t+hx; t+hx, \beta) + o(h^{p+1}). \end{aligned}$$

Thus, using the continuity of $g_2(t; t, \beta)$ and $f_T(t)$ in t , we obtain:

$$\mathbb{E}\{W_n(t, \beta)\} = -c_p \frac{\eta_{\beta}^{(p+1)}(t)}{(p+1)!} g_2(t; t, \beta) f_T(t) h^{p+1} + o(h^{p+1})$$

uniformly in (t, β) . Thus, we conclude that $\sqrt{nh} W_n(t, \beta) = O_p(1)$ when $nh^{2p+3} = O(1)$.

By Lemma A2,

$$\sup_{a^* \in \Theta, t \in \mathcal{T}_0, \beta \in \mathcal{K}} \left| G_n(a^*; t, \beta) - \sqrt{nh} W_n(t, \beta)^T a^* - \frac{1}{2} a^{*T} \mathbf{S}_2(t, \beta) a^* \right| = o_p(1).$$

This along with Lemma A.1 of [18] yields:

$$\sup_{t \in \mathcal{T}_0, \beta \in \mathcal{K}} \|\hat{a}^*(t, \beta) + \{\mathbf{S}_2(t, \beta)\}^{-1} \sqrt{nh} W_n(t, \beta)\| = o_p(1),$$

the first entry of which satisfies:

$$\sup_{t \in \mathcal{T}_0, \beta \in \mathcal{K}} |\sqrt{nh} \{\hat{\eta}_{\beta}(t) - \eta_{\beta}(t)\} + e_{1,p+1}^T \{\mathbf{S}_2(t, \beta)\}^{-1} \sqrt{nh} W_n(t, \beta)| = o_p(1),$$

namely, $\sup_{t \in \mathcal{T}_0, \beta \in \mathcal{K}} |\hat{\eta}_{\beta}(t) - \eta_{\beta}(t) + e_{1,p+1}^T \{\mathbf{S}_2(t, \beta)\}^{-1} W_n(t, \beta)| = o_p(1/\sqrt{nh})$. By [31],

$\sup_{t \in \mathcal{T}_0, \beta \in \mathcal{K}} \|W_n(t, \beta) - \mathbb{E}\{W_n(t, \beta)\}\| = O_p(\{\frac{\log(1/h)}{nh}\}^{1/2})$. Furthermore,

$$\{\mathbf{S}_2(t, \beta)\}^{-1} \mathbb{E}\{W_n(t, \beta)\} = -\mathcal{S}^{-1} c_p \frac{\eta_{\beta}^{(p+1)}(t)}{(p+1)!} h^{p+1} + o(h^{p+1})$$

uniformly in (t, β) . Therefore,

$$\sup_{t \in \mathcal{T}_0, \beta \in \mathcal{K}} \left| \hat{\eta}_{\beta}(t) - \eta_{\beta}(t) - e_{1,p+1}^T \mathcal{S}^{-1} c_p \frac{\eta_{\beta}^{(p+1)}(t)}{(p+1)!} h^{p+1} \right| = o_p(1).$$

This yields:

$$\begin{aligned} & \sup_{\beta \in \mathcal{K}} \sup_{t \in \mathcal{T}_0} \left| \hat{\eta}_{\beta}(t) - \eta_{\beta}(t) - \mathbf{e}_{1,p+1}^T \mathcal{S}^{-1} \mathbf{c}_p \frac{\eta_{\beta}^{(p+1)}(t)}{(p+1)!} h^{p+1} \right| \\ & \leq \sup_{t \in \mathcal{T}_0, \beta \in \mathcal{K}} \left| \hat{\eta}_{\beta}(t) - \eta_{\beta}(t) - \mathbf{e}_{1,p+1}^T \mathcal{S}^{-1} \mathbf{c}_p \frac{\eta_{\beta}^{(p+1)}(t)}{(p+1)!} h^{p+1} \right| = o_p(1). \end{aligned}$$

Note that for $p = 1$, $\mathbf{e}_{1,p+1}^T \mathcal{S}^{-1} \mathbf{c}_p = \mu_2(K)$. This completes the proof. \square

Lemma A1. Assume Condition B in the Appendix. If $n \rightarrow \infty$, $h \rightarrow 0$ and $nh \rightarrow \infty$, then for given $t \in \mathcal{T}_0$ and $\beta \in \mathcal{K}$,

$$\frac{1}{n} \sum_{i=1}^n \mathbf{p}_2(Y_i; \eta_i(t; \beta)) \{ \mathbf{t}_i^*(t) \mathbf{t}_i^*(t)^T \} w_1(\mathbf{X}_i) K_h(T_i - t) = \mathbf{S}_2(t, \beta) + o_p(1),$$

where $\mathbf{S}_2(t, \beta) = g_2(t; t, \beta) f_T(t) \mathcal{S}$, with $\mathcal{S} = (\mu_{j+k-2}(K))_{1 \leq j, k \leq p+1}$ and $\mu_j(K) = \int u^j K(u) du$, $j = 0, 1, \dots, 2p$.

Proof. Recall the $(p+1) \times (p+1)$ matrix $\mathbf{t}_i^*(t) \mathbf{t}_i^*(t)^T = ((\frac{T_i-t}{h})^{j+k-2})_{1 \leq j, k \leq p+1}$. Set $X_j = \frac{1}{n} \sum_{i=1}^n \mathbf{p}_2(Y_i; \eta_i(t; \beta)) (\frac{T_i-t}{h})^j w_1(\mathbf{X}_i) K_h(T_i - t)$ for $j = 0, 1, \dots, 2p$. We observe that:

$$\begin{aligned} E(X_j) &= \frac{1}{n} \sum_{i=1}^n E \left[E \{ \mathbf{p}_2(Y_i; \eta_i(t; \beta)) w_1(\mathbf{X}_i) \mid T_i \} \left(\frac{T_i-t}{h} \right)^j K_h(T_i - t) \right] \\ &= \frac{1}{n} \sum_{i=1}^n E \left\{ g_2(T_i; t, \beta) \left(\frac{T_i-t}{h} \right)^j K_h(T_i - t) \right\} \\ &= E \left\{ g_2(T; t, \beta) \left(\frac{T-t}{h} \right)^j K_h(T - t) \right\} \\ &= \int g_2(y; t, \beta) \left(\frac{y-t}{h} \right)^j \frac{1}{h} K \left(\frac{y-t}{h} \right) f_T(y) dy \\ &= \int g_2(t + hx; t, \beta) x^j K(x) f_T(t + hx) dx \\ &= g_2(t; t, \beta) f_T(t) \mu_j(K) + o(1), \end{aligned}$$

using the continuity of $g_2(\tau; t, \beta)$ in τ and $f_T(t)$ in t . Similarly,

$$\begin{aligned} \text{var}(X_j) &= \frac{1}{n^2} \sum_{i=1}^n \text{var} \left\{ \mathbf{p}_2(Y_i; \eta_i(t; \beta)) \left(\frac{T_i-t}{h} \right)^j w_1(\mathbf{X}_i) K_h(T_i - t) \right\} \\ &\leq \frac{1}{n^2} \sum_{i=1}^n E \left[\mathbf{p}_2^2(Y_i; \eta_i(t; \beta)) \left(\frac{T_i-t}{h} \right)^{2j} w_1^2(\mathbf{X}_i) \{ K_h(T_i - t) \}^2 \right] \\ &\leq \frac{C}{nh}. \end{aligned}$$

This completes the proof. \square

Lemma A2. Assume Condition B. If $n \rightarrow \infty$, $h \rightarrow 0$, $nh \rightarrow \infty$, $\log(1/h)/(nh) \rightarrow 0$, then $\sup_{\mathbf{a}^* \in \Theta, t \in \mathcal{T}_0, \beta \in \mathcal{K}} |G_n(\mathbf{a}^*; t, \beta) - \sqrt{nh} \mathbf{W}_n(t, \beta)^T \mathbf{a}^* - 2^{-1} \mathbf{a}^{*T} \mathbf{S}_2(t, \beta) \mathbf{a}^*| = o_p(1)$, with a compact set $\Theta \subseteq \mathbb{R}^{p+1}$.

Proof. Let $D_n(\mathbf{a}^*; t, \beta) = G_n(\mathbf{a}^*; t, \beta) - \sqrt{nh} \mathbf{W}_n(t, \beta)^T \mathbf{a}^*$. Note that:

$$\begin{aligned} & D_n(\mathbf{a}^*; t, \beta) \\ &= nh \left[\frac{1}{n} \sum_{i=1}^n \rho_q(Y_i, F^{-1}(\eta_i(t; \beta) + \{a_n \mathbf{t}_i^*(t)^T \mathbf{a}^*\})) w_1(\mathbf{X}_i) K_h(T_i - t) \right. \\ & \quad \left. - \frac{1}{n} \sum_{i=1}^n \rho_q(Y_i, F^{-1}(\eta_i(t; \beta))) w_1(\mathbf{X}_i) K_h(T_i - t) \right] \end{aligned}$$

$$\begin{aligned}
& -\frac{1}{n} \sum_{i=1}^n p_1(Y_i; \eta_i(t; \beta)) \{a_n \mathbf{t}_i^*(t)^T \mathbf{a}^*\} w_1(\mathbf{X}_i) K_h(T_i - t) \Big] \\
& = \frac{1}{2} \mathbf{a}^{*T} \left[\frac{1}{n} \sum_{i=1}^n p_2(Y_i; \tilde{\eta}_i(t; \beta)) \{\mathbf{t}_i^*(t) \mathbf{t}_i^*(t)^T\} w_1(\mathbf{X}_i) K_h(T_i - t) \right] \mathbf{a}^*,
\end{aligned}$$

where $a_n = 1/\sqrt{nh}$ and $\tilde{\eta}_i(t; \beta)$ is between $\eta_i(t; \beta)$ and $\eta_i(t; \beta) + \{a_n \mathbf{t}_i^*(t)^T \mathbf{a}^*\}$. Then:

$$\begin{aligned}
& |D_n(\mathbf{a}^*; t, \beta) - 2^{-1} \mathbf{a}^{*T} \mathbf{S}_2(t, \beta) \mathbf{a}^*| \\
& = \frac{1}{2} \left| \mathbf{a}^{*T} \left[\frac{1}{n} \sum_{i=1}^n p_2(Y_i; \tilde{\eta}_i(t; \beta)) \{\mathbf{t}_i^*(t) \mathbf{t}_i^*(t)^T\} w_1(\mathbf{X}_i) K_h(T_i - t) - \mathbf{S}_2(t, \beta) \right] \mathbf{a}^* \right| \\
& \leq \|\mathbf{a}^*\|^2 \left| \frac{1}{n} \sum_{i=1}^n p_2(Y_i; \tilde{\eta}_i(t; \beta)) \{\mathbf{t}_i^*(t) \mathbf{t}_i^*(t)^T\} w_1(\mathbf{X}_i) K_h(T_i - t) - \mathbf{S}_2(t, \beta) \right|.
\end{aligned}$$

The proof completes by applying [31]. \square

Proof of Theorem 1. Before showing Theorem 1, we need Proposition A1 (whose proof is omitted), where the following notation will be used. Denote by $\mathcal{C}^1(\mathcal{T})$ the set of continuously differentiable functions in \mathcal{T} . Let $\mathcal{V}(\beta)$ denote the neighborhood of $\beta \in \mathcal{K}$. Let $\mathcal{H}_\delta(\beta)$ denote the neighborhood of η_β such that $\mathcal{V}(\beta) \subseteq \mathcal{K}$ and $\mathcal{H}_\delta(\beta) = \{u \in \mathcal{C}^1(\mathcal{T}) : \|u - \eta_\beta\|_\infty \leq \delta, \|\frac{\partial}{\partial t} u - \frac{\partial}{\partial t} \eta_\beta\|_\infty \leq \delta\}$.

Proposition A1. Let $\{(Y_i, \mathbf{X}_i, T_i)\}_{i=1}^n$ be independent observations of (Y, \mathbf{X}, T) modeled by (2) and (5). Assume that a random variable T is distributed on \mathcal{T} . Let \mathcal{K} and $\mathcal{H}_1(\beta)$ be compact sets, $g(\cdot; \cdot) : \mathbb{R}^2 \rightarrow \mathbb{R}$ be a continuous and bounded function, $W(\mathbf{x}, t) : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$ be such that $E\{|W(\mathbf{X}, T)|\} < \infty$ and $\eta_\beta(t) = \eta(t, \beta) : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$ be a continuous function of (t, β) . Then:

- (i) $E\{g(Y; \mathbf{X}^T \boldsymbol{\theta} + v(T)) W(\mathbf{X}, T)\} \rightarrow E\{g(Y; \mathbf{X}^T \boldsymbol{\beta} + \eta_\beta(T)) W(\mathbf{X}, T)\}$ as $\|\boldsymbol{\theta} - \boldsymbol{\beta}\| + \|v - \eta_\beta\|_\infty \rightarrow 0$;
- (ii) $\sup_{\boldsymbol{\theta} \in \mathcal{K}} |n^{-1} \sum_{i=1}^n g(Y_i; \mathbf{X}_i^T \boldsymbol{\theta} + \eta_\theta(T_i)) W(\mathbf{X}_i, T) - E\{g(Y; \mathbf{X}^T \boldsymbol{\theta} + \eta_\theta(T)) W(\mathbf{X}, T)\}| \xrightarrow{P} 0$ as $n \rightarrow \infty$;
- (iii) if, in addition, \mathcal{T} is compact and $\eta_\beta \in \mathcal{C}^1(\mathcal{T})$, then $\sup_{\boldsymbol{\theta} \in \mathcal{K}, v \in \mathcal{H}_1(\beta)} |n^{-1} \sum_{i=1}^n g(Y_i; \mathbf{X}_i^T \boldsymbol{\theta} + v(T_i)) W(\mathbf{X}_i, T_i) - E\{g(Y; \mathbf{X}^T \boldsymbol{\theta} + v(T)) W(\mathbf{X}, T)\}| \xrightarrow{P} 0$ as $n \rightarrow \infty$.

For part (i), we first show that for any compact set \mathcal{K} in \mathbb{R}^d ,

$$\sup_{\beta \in \mathcal{K}} |J_n(\beta, \hat{\eta}_\beta) - J(\beta, \eta_\beta)| \xrightarrow{P} 0. \quad (\text{A3})$$

It suffices to show $\sup_{\beta \in \mathcal{K}} |J_n(\beta, \eta_\beta) - J(\beta, \eta_\beta)| \xrightarrow{P} 0$, which follows from Proposition A1 (ii), and:

$$\sup_{\beta \in \mathcal{K}} |J_n(\beta, \hat{\eta}_\beta) - J_n(\beta, \eta_\beta)| \xrightarrow{P} 0. \quad (\text{A4})$$

To show (A4), we note that for any $\varepsilon > 0$, let \mathcal{T}_0 be a compact set such that $P(T_i \notin \mathcal{T}_0) < \varepsilon$. Then:

$$\begin{aligned}
& J_n(\beta, \hat{\eta}_\beta) - J_n(\beta, \eta_\beta) \\
& = \frac{1}{n} \sum_{i=1}^n \{\rho_q(Y_i, F^{-1}(\mathbf{X}_i^T \boldsymbol{\beta} + \hat{\eta}_\beta(T_i))) - \rho_q(Y_i, F^{-1}(\mathbf{X}_i^T \boldsymbol{\beta} + \eta_\beta(T_i)))\} w_2(\mathbf{X}_i) \mathbf{I}(T_i \in \mathcal{T}_0) \\
& \quad + \frac{1}{n} \sum_{i=1}^n \{\rho_q(Y_i, F^{-1}(\mathbf{X}_i^T \boldsymbol{\beta} + \hat{\eta}_\beta(T_i))) - \rho_q(Y_i, F^{-1}(\mathbf{X}_i^T \boldsymbol{\beta} + \eta_\beta(T_i)))\} w_2(\mathbf{X}_i) \mathbf{I}(T_i \notin \mathcal{T}_0).
\end{aligned}$$

For $T_i \in \mathcal{T}_0$, by the mean-value theorem,

$$\begin{aligned}
& |\rho_q(Y_i, F^{-1}(\mathbf{X}_i^T \boldsymbol{\beta} + \hat{\eta}_\beta(T_i))) - \rho_q(Y_i, F^{-1}(\mathbf{X}_i^T \boldsymbol{\beta} + \eta_\beta(T_i)))| \\
& = |p_1(Y_i; \mathbf{X}_i^T \boldsymbol{\beta} + \eta_{i,\beta}^*) \{\hat{\eta}_\beta(T_i) - \eta_\beta(T_i)\}|
\end{aligned}$$

$$\leq \|p_1(\cdot; \cdot)\|_\infty \sup_{\beta \in \mathcal{K}} \|\hat{\eta}_\beta - \eta_\beta\|_{\mathcal{T}_0; \infty},$$

where $\eta_{i,\beta}^*$ is located between $\hat{\eta}_\beta(T_i)$ and $\eta_\beta(T_i)$. For $T_i \notin \mathcal{T}_0$, it follows that:

$$\begin{aligned} & |\rho_q(Y_i, F^{-1}(X_i^T \beta + \hat{\eta}_\beta(T_i))) - \rho_q(Y_i, F^{-1}(X_i^T \beta + \eta_\beta(T_i)))| \\ & \leq 2\|\rho_q(\cdot, \cdot)\|_\infty. \end{aligned}$$

Hence,

$$\begin{aligned} |J_n(\beta, \hat{\eta}_\beta) - J_n(\beta, \eta_\beta)| & \leq \left\{ \|p_1(\cdot; \cdot)\|_\infty \sup_{\beta \in \mathcal{K}} \|\hat{\eta}_\beta - \eta_\beta\|_{\mathcal{T}_0; \infty} + 2\|\rho_q(\cdot, \cdot)\|_\infty T_n^* \right\} \|w_2\|_\infty \\ & \leq 2\varepsilon, \end{aligned}$$

where the last inequality is entailed by Lemma 1 and the law of large numbers for $T_n^* = n^{-1} \sum_{i=1}^n \mathbf{I}(T_i \notin \mathcal{T}_0)$. This completes the proof of (A3). The proof of $\hat{\beta} \xrightarrow{P} \beta_0$ follows from combining Lemma A-1 of [1] with (A3) and Condition A2.

Part (ii) follows from Lemma 1, Part (i) and Condition B5 for $\eta_\beta(t)$. \square

Proof of Theorem 2. Similar to the proof of Lemma 1, it can be shown that $|\hat{\eta}_\beta(t) - \eta_\beta(t) + e_{1,p+1}^T \{S_2(t, \beta)\}^{-1} \frac{1}{n} \sum_{i=1}^n p_1(Y_i; \eta_i(t; \beta)) \mathbf{t}_i^*(t) w_1(X_i) K_h(T_i - t)| = O_p(h^2 a_n + a_n^2 \sqrt{\log(1/h)})$. Note that for $p = 1$,

$$\begin{aligned} e_{1,p+1}^T \{S_2(t, \beta)\}^{-1} \mathbf{t}_i^*(t) & = \frac{1}{g_2(t; t, \beta) f_T(t)} (1, 0) \begin{pmatrix} 1 & 0 \\ 0 & 1/\mu_2(K) \end{pmatrix} \begin{pmatrix} 1 \\ (T_i - t)/h \end{pmatrix} \\ & = \frac{1}{g_2(t; t, \beta) f_T(t)}. \end{aligned}$$

Thus:

$$\left| \hat{\eta}_\beta(t) - \eta_\beta(t) + \frac{1}{n f_T(t) g_2(t; t, \beta)} \sum_{i=1}^n p_1(Y_i; \eta_i(t; \beta)) w_1(X_i) K_h(T_i - t) \right| = O_p(h^2 a_n + a_n^2 \sqrt{\log(1/h)}).$$

Consider $\hat{\beta}$ defined in (23). Note that:

$$\begin{aligned} X_i^T \beta + \hat{\eta}_\beta(T_i) & = X_i^T \beta_0 + X_i^T (\beta - \beta_0) + \hat{\eta}_{(\beta - \beta_0) + \beta_0}(T_i) \\ & = X_i^T \beta_0 + c_n X_i^T \{\sqrt{n}(\beta - \beta_0)\} + \hat{\eta}_{c_n \{\sqrt{n}(\beta - \beta_0)\} + \beta_0}(T_i), \end{aligned}$$

where $c_n = 1/\sqrt{n}$. Then, $\hat{\theta} = \sqrt{n}(\hat{\beta} - \beta_0)$ minimizes:

$$\begin{aligned} J_n(\theta) & = n \left[\frac{1}{n} \sum_{i=1}^n \left\{ \rho_q(Y_i, F^{-1}(X_i^T \beta_0 + c_n X_i^T \theta + \hat{\eta}_{c_n \theta + \beta_0}(T_i))) w_2(X_i) \right. \right. \\ & \quad \left. \left. - \rho_q(Y_i, F^{-1}(X_i^T \beta_0 + \hat{\eta}_{\beta_0}(T_i))) w_2(X_i) \right\} \right] \end{aligned}$$

with respect to θ . By Taylor expansion,

$$\begin{aligned} & J_n(\theta) \\ & = n \left(\frac{1}{n} \sum_{i=1}^n p_1(Y_i; X_i^T \beta_0 + \hat{\eta}_{\beta_0}(T_i)) [c_n X_i^T \theta + \{\hat{\eta}_{c_n \theta + \beta_0}(T_i) - \hat{\eta}_{\beta_0}(T_i)\}] w_2(X_i) \right) \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{2n} \sum_{i=1}^n p_2(Y_i; \mathbf{X}_i^T \boldsymbol{\beta}_0 + \hat{\eta}_{\beta_0}(T_i)) [c_n \mathbf{X}_i^T \boldsymbol{\theta} + \{\hat{\eta}_{c_n \boldsymbol{\theta} + \beta_0}(T_i) - \hat{\eta}_{\beta_0}(T_i)\}]^2 w_2(\mathbf{X}_i) \\
& + \frac{1}{6n} \sum_{i=1}^n p_3(Y_i; \eta_i^*) [c_n \mathbf{X}_i^T \boldsymbol{\theta} + \{\hat{\eta}_{c_n \boldsymbol{\theta} + \beta_0}(T_i) - \hat{\eta}_{\beta_0}(T_i)\}]^3 w_2(\mathbf{X}_i) \Big) \\
& = I_{n,1} + I_{n,2} + I_{n,3},
\end{aligned}$$

where η_i^* is located between $\mathbf{X}_i^T \boldsymbol{\beta}_0 + \hat{\eta}_{\beta_0}(T_i)$ and $\mathbf{X}_i^T \boldsymbol{\beta}_0 + c_n \mathbf{X}_i^T \boldsymbol{\theta} + \hat{\eta}_{c_n \boldsymbol{\theta} + \beta_0}(T_i)$,

$$\begin{aligned}
I_{n,1} &= \sum_{i=1}^n p_1(Y_i; \mathbf{X}_i^T \boldsymbol{\beta}_0 + \hat{\eta}_{\beta_0}(T_i)) \left\{ c_n \mathbf{X}_i^T \boldsymbol{\theta} + \frac{\partial \hat{\eta}_{\beta}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_n} c_n \boldsymbol{\theta} \right\} w_2(\mathbf{X}_i) \\
&= \frac{1}{\sqrt{n}} \sum_{i=1}^n p_1(Y_i; \mathbf{X}_i^T \boldsymbol{\beta}_0 + \hat{\eta}_{\beta_0}(T_i)) \left\{ \mathbf{X}_i + \frac{\partial \hat{\eta}_{\beta}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_n} \right\}^T \boldsymbol{\theta} w_2(\mathbf{X}_i) \\
&= \frac{1}{\sqrt{n}} \sum_{i=1}^n p_1(Y_i; \mathbf{X}_i^T \boldsymbol{\beta}_0 + \hat{\eta}_{\beta_0}(T_i)) \left\{ \mathbf{X}_i + \frac{\partial \eta_{\beta}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\}^T \boldsymbol{\theta} w_2(\mathbf{X}_i) + o_p(1), \\
I_{n,2} &= \frac{1}{2} \boldsymbol{\theta}^T \left[\frac{1}{n} \sum_{i=1}^n p_2(Y_i; \mathbf{X}_i^T \boldsymbol{\beta}_0 + \hat{\eta}_{\beta_0}(T_i)) \right. \\
&\quad \left. \left\{ \mathbf{X}_i + \frac{\partial \hat{\eta}_{\beta}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_n} \right\} \left\{ \mathbf{X}_i + \frac{\partial \hat{\eta}_{\beta}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_n} \right\}^T w_2(\mathbf{X}_i) \right] \boldsymbol{\theta} \\
&= \frac{1}{2} \boldsymbol{\theta}^T \mathbf{B}_2 \boldsymbol{\theta} + o_p(1), \\
I_{n,3} &= o_p(1),
\end{aligned}$$

with $\boldsymbol{\beta}_n$ located between $\boldsymbol{\beta}_0$ and $c_n \boldsymbol{\theta} + \boldsymbol{\beta}_0$, and $\mathbf{B}_2 = \mathbf{H}_0$ following Lemma 1, Condition A3 and Proposition A1. Thus:

$$J_n(\boldsymbol{\theta}) = I_{n,1}^* \boldsymbol{\theta} + \frac{1}{2} \boldsymbol{\theta}^T \mathbf{B}_2 \boldsymbol{\theta} + o_p(1), \quad (\text{A5})$$

where $I_{n,1}^* = \frac{1}{\sqrt{n}} \sum_{i=1}^n p_1(Y_i; \mathbf{X}_i^T \boldsymbol{\beta}_0 + \hat{\eta}_{\beta_0}(T_i)) \left\{ \mathbf{X}_i + \frac{\partial \eta_{\beta}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\} w_2(\mathbf{X}_i)$. Note that:

$$\begin{aligned}
I_{n,1}^* &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[p_1(Y_i; \mathbf{X}_i^T \boldsymbol{\beta}_0 + \eta_{\beta_0}(T_i)) \left\{ \mathbf{X}_i + \frac{\partial \eta_{\beta}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\} w_2(\mathbf{X}_i) \right. \\
&\quad + p_2(Y_i; \mathbf{X}_i^T \boldsymbol{\beta}_0 + \eta_{\beta_0}(T_i)) \left\{ \mathbf{X}_i + \frac{\partial \eta_{\beta}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\} w_2(\mathbf{X}_i) \{\hat{\eta}_{\beta_0}(T_i) - \eta_{\beta_0}(T_i)\} \\
&\quad \left. + \frac{1}{2} p_3(Y_i; \eta_i^{**}) \left\{ \mathbf{X}_i + \frac{\partial \eta_{\beta}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\} w_2(\mathbf{X}_i) \{\hat{\eta}_{\beta_0}(T_i) - \eta_{\beta_0}(T_i)\}^2 \right] \\
&= T_{n,1} + T_{n,2} + T_{n,3},
\end{aligned}$$

where η_i^{**} is between $\mathbf{X}_i^T \boldsymbol{\beta}_0 + \hat{\eta}_{\beta_0}(T_i)$ and $\mathbf{X}_i^T \boldsymbol{\beta}_0 + \eta_{\beta_0}(T_i)$,

$$\begin{aligned}
T_{n,3} &= o_p(1), \\
T_{n,2} &= \frac{1}{\sqrt{n}} \sum_{i=1}^n p_2(Y_i; \mathbf{X}_i^T \boldsymbol{\beta}_0 + \eta_{\beta_0}(T_i)) \left\{ \mathbf{X}_i + \frac{\partial \eta_{\beta}(T_i)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\} w_2(\mathbf{X}_i) \\
&\quad \times \frac{(-1)}{n f_T(T_i) g_2(T_i; T_i, \boldsymbol{\beta}_0)} \sum_{j=1}^n p_1(Y_j; \eta_j(T_i; \boldsymbol{\beta}_0)) w_1(\mathbf{X}_j) K_h(T_j - T_i) \\
&= -\frac{1}{\sqrt{n}} \sum_{j=1}^n \frac{p_1(Y_j; \eta_j) w_1(\mathbf{X}_j)}{g_2(T_j; T_j, \boldsymbol{\beta}_0)} E \left[p_2(Y_j; \eta_j) \left\{ \mathbf{X}_j + \frac{\partial \eta_{\beta}(T_j)}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \right\} w_2(\mathbf{X}_j) \Big| T_j \right] \\
&\equiv -\frac{1}{\sqrt{n}} \sum_{j=1}^n p_1(Y_j; \eta_j) \frac{\gamma(T_j)}{g_2(T_j; T_j, \boldsymbol{\beta}_0)} w_1(\mathbf{X}_j),
\end{aligned}$$

with:

$$\gamma(t) = E \left[p_2(Y; \eta(X, T)) \left\{ X + \frac{\partial \eta_\beta(T)}{\partial \beta} \Big|_{\beta=\beta_0} \right\} w_2(X) \Big| T = t \right].$$

Therefore,

$$I_{n,1}^* = \frac{1}{\sqrt{n}} \sum_{i=1}^n p_1(Y_i; \eta_i) \left[\left\{ X_i + \frac{\partial \eta_\beta(T_i)}{\partial \beta} \Big|_{\beta=\beta_0} \right\} w_2(X_i) - \frac{\gamma(T_i)}{g_2(T_i; T_i, \beta_0)} w_1(X_i) \right] + o_p(1).$$

By the central limit theorem,

$$I_{n,1}^* \xrightarrow{\mathcal{D}} N(\mathbf{0}, \Omega_0^*), \quad (\text{A6})$$

where:

$$\begin{aligned} \Omega_0^* &= E \left(p_1^2(Y; \eta(X, T)) \left[\left\{ X + \frac{\partial \eta_\beta(T)}{\partial \beta} \Big|_{\beta=\beta_0} \right\} w_2(X) - \frac{\gamma(T)}{g_2(T; T, \beta_0)} w_1(X) \right] \right. \\ &\quad \left. \left[\left\{ X + \frac{\partial \eta_\beta(T)}{\partial \beta} \Big|_{\beta=\beta_0} \right\} w_2(X) - \frac{\gamma(T)}{g_2(T; T, \beta_0)} w_1(X) \right]^T \right). \end{aligned}$$

From (A5) and (A6), $\hat{\theta} = -\mathbf{B}_2^{-1} I_{n,1}^* + o_p(1)$. This implies that $\sqrt{n}(\hat{\beta} - \beta_0) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{H}_0^{-1} \Omega_0^* \mathbf{H}_0^{-1})$. \square

Proof of Theorem 3. Denote $\mathbf{V}_0 = \mathbf{H}_0^{-1} \Omega_0^* \mathbf{H}_0^{-1}$ and $\hat{\mathbf{V}}_n = \hat{\mathbf{H}}_0^{-1} \hat{\Omega}_0^* \hat{\mathbf{H}}_0^{-1}$. Note that $\mathbf{A}\hat{\beta} - \mathbf{g}_0 = \mathbf{A}(\hat{\beta} - \beta_0) + (\mathbf{A}\beta_0 - \mathbf{g}_0)$. Thus:

$$\begin{aligned} & (\mathbf{A}\hat{\mathbf{V}}_n \mathbf{A}^T)^{-1/2} \sqrt{n}(\mathbf{A}\hat{\beta} - \mathbf{g}_0) \\ &= (\mathbf{A}\hat{\mathbf{V}}_n \mathbf{A}^T)^{-1/2} \{ \mathbf{A}\sqrt{n}(\hat{\beta} - \beta_0) \} + (\mathbf{A}\hat{\mathbf{V}}_n \mathbf{A}^T)^{-1/2} \{ \sqrt{n}(\mathbf{A}\beta_0 - \mathbf{g}_0) \} \\ &\equiv I_1 + I_2, \end{aligned}$$

which implies that $W_n = \|I_1 + I_2\|^2$. Arguments for Theorem 2 give $I_1 \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{I}_k)$. Under H_0 in (4), $I_2 \equiv \mathbf{0}$ and thus $(I_1 + I_2) \xrightarrow{\mathcal{D}} N(\mathbf{0}, \mathbf{I}_k)$, which completes the proof. \square

Proof of Theorem 4. Follow the notation and proof in Theorem 3. Under H_{1n} in (29), $I_2 \xrightarrow{\mathcal{P}} (\mathbf{A}\mathbf{V}_0 \mathbf{A}^T)^{-1/2} \mathbf{c}$ and thus $(I_1 + I_2) \xrightarrow{\mathcal{D}} N((\mathbf{A}\mathbf{V}_0 \mathbf{A}^T)^{-1/2} \mathbf{c}, \mathbf{I}_k)$. This completes the proof. \square

Proof of Theorem 5. Following the notation and proof in Theorem 3, $W_n = \|I_1\|^2 + 2I_1^T I_2 + \|I_2\|^2$. We see that $\|I_1\|^2 \xrightarrow{\mathcal{D}} \chi_k^2$. Under H_1 in (4), $I_2 = (\mathbf{A}\mathbf{V}_0 \mathbf{A}^T)^{-1/2} \sqrt{n}(\mathbf{A}\beta_0 - \mathbf{g}_0) \{1 + o_p(1)\}$, which means $\|I_2\|^2 = n(\mathbf{A}\beta_0 - \mathbf{g}_0)^T (\mathbf{A}\mathbf{V}_0 \mathbf{A}^T)^{-1} (\mathbf{A}\beta_0 - \mathbf{g}_0) \{1 + o_p(1)\}$ and thus $I_1^T I_2 = O_p(\sqrt{n})$. Hence, $n^{-1}W_n \geq \lambda_{\min}\{(\mathbf{A}\mathbf{V}_0 \mathbf{A}^T)^{-1}\} \|\mathbf{A}\beta_0 - \mathbf{g}_0\|^2 + o_p(1)$. This completes the proof. \square

Proof of Theorem 6. Denote $J_n(\beta) = J_n(\beta, \hat{\eta}_\beta)$. For the matrix \mathbf{A} in (4), there exists a $(d-k) \times d$ matrix \mathbf{B} satisfying $\mathbf{B}\mathbf{B}^T = \mathbf{I}_{d-k}$ and $\mathbf{A}\mathbf{B}^T = \mathbf{0}$. Therefore, $\mathbf{A}\beta = \mathbf{g}_0$ is equivalent to $\beta = \mathbf{B}^T \gamma + \mathbf{b}_0$ for some vector $\gamma \in \mathbb{R}^{d-k}$ and $\mathbf{b}_0 = \mathbf{A}^T (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{g}_0$. Then, minimizing $J_n(\beta)$ subject to $\mathbf{A}\beta = \mathbf{g}_0$ is equivalent to minimizing $J_n(\mathbf{B}^T \gamma + \mathbf{b}_0)$ with respect to γ , and we denote by $\hat{\gamma}$ the minimizer. Furthermore, under H_0 in (4), we have $\beta_0 = \mathbf{B}^T \gamma_0 + \mathbf{b}_0$ for $\gamma_0 = \mathbf{B}\beta_0$, and $\hat{\gamma} - \gamma_0 \xrightarrow{\mathcal{P}} \mathbf{0}$.

For Part (i), using the Taylor expansion around $\hat{\beta}$, we get:

$$J_n(\mathbf{B}^T \hat{\gamma} + \mathbf{b}_0) - J_n(\hat{\beta}) = \frac{1}{2n} \{ \sqrt{n}(\mathbf{B}^T \hat{\gamma} + \mathbf{b}_0 - \hat{\beta}) \}^T J_n''(\tilde{\beta}) \{ \sqrt{n}(\mathbf{B}^T \hat{\gamma} + \mathbf{b}_0 - \hat{\beta}) \}, \quad (\text{A7})$$

where $\tilde{\beta}$ is between $B^T \hat{\gamma} + \mathbf{b}_0$ and $\hat{\beta}$. We now discuss $B^T \hat{\gamma} + \mathbf{b}_0 - \hat{\beta}$. From the proof in Theorem 2, $(\hat{\beta} - \beta_0) = -\mathbf{H}_0^{-1} J'_n(\beta_0) \{1 + o_p(1)\}$, where $J'_n(\beta_0) = \{J'_{n,1} + o_p(1)\} / \sqrt{n}$. Similar arguments deduce $\hat{\gamma} - \gamma_0 = -(B\mathbf{H}_0 B^T)^{-1} B J'_n(\beta_0) \{1 + o_p(1)\}$. Thus, under H_0 in (4),

$$B^T \hat{\gamma} + \mathbf{b}_0 - \hat{\beta} = B^T (\hat{\gamma} - \gamma_0) - (\hat{\beta} - \beta_0) = \mathbf{H}_0^{-1/2} P_{\mathbf{H}_0^{-1/2} \mathbf{A}^T} \mathbf{H}_0^{-1/2} J'_n(\beta_0) \{1 + o_p(1)\},$$

and thus by (A6),

$$\sqrt{n}(B^T \hat{\gamma} + \mathbf{b}_0 - \hat{\beta}) \xrightarrow{\mathcal{D}} \mathbf{H}_0^{-1/2} P_{\mathbf{H}_0^{-1/2} \mathbf{A}^T} \mathbf{H}_0^{-1/2} \Omega_0^{*1/2} \mathbf{Z}, \quad (\text{A8})$$

where $\mathbf{Z} = (Z_1, \dots, Z_d)^T \sim N(\mathbf{0}, \mathbf{I}_d)$. Combining the fact $J''_n(\tilde{\beta}) \xrightarrow{\text{P}} \mathbf{H}_0$, (A7) and (A8) gives:

$$\begin{aligned} \Lambda_n &= \{\sqrt{n}(B^T \hat{\gamma} + \mathbf{b}_0 - \hat{\beta})\}^T \mathbf{H}_0 \{\sqrt{n}(B^T \hat{\gamma} + \mathbf{b}_0 - \hat{\beta})\} \{1 + o_p(1)\} \\ &\xrightarrow{\mathcal{D}} \mathbf{Z}^T \Omega_0^{*1/2} \mathbf{H}_0^{-1/2} P_{\mathbf{H}_0^{-1/2} \mathbf{A}^T} \mathbf{H}_0^{-1/2} \Omega_0^{*1/2} \mathbf{Z} \\ &= \sum_{j=1}^d \lambda_j (\Omega_0^{*1/2} \mathbf{H}_0^{-1/2} P_{\mathbf{H}_0^{-1/2} \mathbf{A}^T} \mathbf{H}_0^{-1/2} \Omega_0^{*1/2}) Z_j^2 \\ &= \sum_{j=1}^k \lambda_j \{(\mathbf{A} \mathbf{H}_0^{-1} \mathbf{A}^T)^{-1} (\mathbf{A} \mathbf{V}_0 \mathbf{A}^T)\} Z_j^2. \end{aligned} \quad (\text{A9})$$

This proves Part (i).

For Part (ii), using $\psi(r) = r$, $w_1(\mathbf{x}) = w_2(\mathbf{x}) \equiv 1$ and (31), we obtain $\Omega_0^* = \Omega_0 = \mathbf{C} \mathbf{H}_0$, and thus, $\mathbf{A} \mathbf{V}_0 \mathbf{A}^T = \mathbf{C} (\mathbf{A} \mathbf{H}_0^{-1} \mathbf{A}^T)$. Thus, (A9) = $\mathbf{C} \sum_{j=1}^k Z_j^2 \sim \mathbf{C} \chi_k^2$, which completes the proof. \square

Proof of Theorem 7. The proofs are similar to those used in Theorem 4 and Theorems 5 and 6. The lengthy details are omitted. \square

References

- Andrews, D. Asymptotics for semiparametric econometric models via stochastic equicontinuity. *Econometrica* **1994**, *62*, 43–72.
- Robinson, P.M. Root-n consistent semiparametric regression. *Econometrica* **1988**, *56*, 931–954.
- Speckman, P. Kernel smoothing in partial linear models. *J. R. Statist. Soc. B* **1988**, *50*, 413–436.
- Yatchew, A. An elementary estimator of the partial linear model. *Econ. Lett.* **1997**, *57*, 135–143.
- Fan, J.; Li, R. New estimation and model selection procedures for semiparametric modeling in longitudinal data analysis. *J. Am. Stat. Assoc.* **2004**, *99*, 710–723.
- McCullagh, P.; Nelder, J.A. *Generalized Linear Models*, 2nd ed.; Chapman & Hall: London, UK, 1989.
- Zhang, C.M.; Yu, T. Semiparametric detection of significant activation for brain fMRI. *Ann. Stat.* **2008**, *36*, 1693–1725.
- Fan, J.; Huang, T. Profile likelihood inferences on semiparametric varying-coefficient partially linear models. *Bernoulli* **2005**, *11*, 1031–1057.
- Boente, G.; He, X.; Zhou, J. Robust estimates in generalized partially linear models. *Ann. Stat.* **2006**, *34*, 2856–2878.
- Zhang, C.M.; Guo, X.; Cheng, C.; Zhang, Z.J. Robust-BD estimation and inference for varying-dimensional general linear models. *Stat. Sin.* **2014**, *24*, 653–673.
- Fan, J.; Gijbels, I. *Local Polynomial Modeling and Its Applications*; Chapman and Hall: London, UK, 1996.
- Brègman, L.M. A relaxation method of finding a common point of convex sets and its application to the solution of problems in convex programming. *USSR Comput. Math. Math. Phys.* **1967**, *7*, 620–631.
- Hastie, T.; Tibshirani, R.; Friedman, J. *The Elements of Statistical Learning*; Springer: New York, NY, USA, 2001.
- Zhang, C.M.; Jiang, Y.; Shang, Z. New aspects of Bregman divergence in regression and classification with parametric and nonparametric estimation. *Can. J. Stat.* **2009**, *37*, 119–139.
- Huber, P. Robust estimation of a location parameter. *Ann. Math. Statist.* **1964**, *35*, 73–101.
- Van der Vaart, A.W. *Asymptotic Statistics*; Cambridge University Press: Cambridge, UK, 1998.

17. Friedman, J.; Hastie, T.; Tibshirani, R. Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.* **2010**, *33*, 1–22.
18. Carroll, R.; Fan, J.; Gijbels, I.; Wand, M. Generalized partially linear single-index models. *J. Am. Stat. Assoc.* **1997**, *92*, 477–489.
19. Mukarjee, H.; Stern, S. Feasible nonparametric estimation of multiargument monotone functions. *J. Am. Stat. Assoc.* **1994**, *89*, 77–80.
20. Albright, S.C.; Winston, W.L.; Zappe, C.J. *Data Analysis and Decision Making with Microsoft Excel*; Duxbury Press: Pacific Grove, CA, USA, 1999.
21. Fan, J.; Peng, H. Nonconcave penalized likelihood with a diverging number of parameters. *Ann. Stat.* **2004**, *32*, 928–961.
22. Dette, H. A consistent test for the functional form of a regression based on a difference of variance estimators. *Ann. Stat.* **1999**, *27*, 1012–1050.
23. Dette, H.; von Lieres und Wilkau, C. Testing additivity by kernel-based methods. *Bernoulli* **2001**, *7*, 669–697.
24. Fan, J.; Zhang, C.M.; Zhang, J. Generalized likelihood ratio statistics and Wilks phenomenon. *Ann. Stat.* **2001**, *29*, 153–193.
25. Hong, Y.M.; Lee, Y.J. A loss function approach to model specification testing and its relative efficiency. *Ann. Stat.* **2013**, *41*, 1166–1203.
26. Zheng, J.X. A consistent test of functional form via nonparametric estimation techniques. *J. Econ.* **1996**, *75*, 263–289.
27. Opsomer, J.D.; Ruppert, D. A root-n consistent backfitting estimator for semiparametric additive modeling. *J. Comput. Graph. Stat.* **1999**, *8*, 715–732.
28. Belloni, A.; Chernozhukov, V.; Hansen, C. Inference on treatment effects after selection amongst high-dimensional controls. *Rev. Econ. Stud.* **2014**, *81*, 608–650.
29. Cattaneo, M.D.; Jansson, M.; Newey, W.K. Alternative asymptotics and the partially linear model with many regressors. *Econ. Theory* **2016**, 1–25.
30. Cattaneo, M.D.; Jansson, M.; Newey, W.K. Treatment effects with many covariates and heteroskedasticity. *arXiv* **2015**, arXiv:1507.02493.
31. Mack, Y.P.; Silverman, B.W. Weak and strong uniform consistency of kernel regression estimates. *Z. Wahrsch. Verw. Gebiete* **1982**, *61*, 405–415.



© 2017 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).