

Adaptively weighted group Lasso for semiparametric quantile regression models

Toshio Honda, Ching-Kang Ing, and Wei-Ying Wu

Abstract

We propose an adaptively weighted group Lasso procedure for simultaneous variable selection and structure identification for varying coefficient quantile regression models and additive quantile regression models with ultra-high dimensional covariates. Under a strong sparsity condition, we establish selection consistency of the proposed Lasso procedure when the weights therein satisfy a set of general conditions. This consistency result, however, is reliant on a suitable choice of the tuning parameter for the Lasso penalty, which can be hard to make in practice. To alleviate this difficulty, we suggest a BIC-type criterion, which we call high-dimensional information criterion (HDIC), and show that the proposed Lasso procedure with the tuning parameter determined by HDIC still achieves selection consistency. Our simulation studies support strongly our theoretical findings.

Keywords: Additive models; B-spline; high-dimensional information criteria; Lasso; structure identification; varying coefficient models.

1 Introduction

We propose adaptively weighted group Lasso (AWG-Lasso) procedures for simultaneous variable selection and structure identification for varying coefficient quantile regression models and additive quantile regression models with ultra-high dimension covariates. Let the number of covariates be denoted by p . Throughout this paper, we assume $p = O(\exp(n^\iota))$, where n is the sample size and ι is a positive constant specified later in Assumption A4 and A4' of Section 5. Under a strong sparsity condition, we establish selection consistency of AWG-Lasso when its weights, determined by some initial estimates, e.g., Lasso and group Lasso, obey a set of general conditions. This consistency

Toshio Honda is Professor, Graduate School of Economics, Hitotsubashi University, 2-1 Naka, Kunitachi, Tokyo 186-8601, Japan (Email: t.honda@r.hit-u.ac.jp). Ching-Kang Ing is Professor, Institute of Statistics, National Tsing Hua University, 101, Sec. 2, Kuang-Fu Rd., Hsinchu 30013, Taiwan (Email: cking@stat.nthu.edu.tw). Wei-Ying Wu is Assistant Professor, Department of Applied Mathematics, National Dong Hwa University, Taiwan (wuweiyong1011@gmail.com). Honda's research was supported by JSPS KAKENHI Grant Number JP 16K05268. Ing's research was supported in part by the Science Vanguard Research Program of the Ministry of Science and Technology (MOST), Taiwan.

result, however, is reliant on a suitable choice for the tuning parameter for the Lasso penalty, which can be hard to make in practice. To alleviate this difficulty, we suggest a BIC-type criterion, which we call high-dimensional information criterion (HDIC), and show that AWG-Lasso with the penalty determined by HDIC (denoted by AWG-Lasso+HDIC hereafter) still achieves selection consistency. This latter result improves previous ones in [21] and the BIC results in [38] since the former does not deal with semi-parametric models and the latter concentrates on linear models. See also [5] and [19] for recent developments in BIC-type model selection criteria. With the selected model, one can conduct final statistical inference by appealing to the results as in [34], [4], [28]. Moreover, our approach can be implemented at several different quantiles, thereby leading to a deeper understanding of the data in hand. There are some other approaches to quantile estimation from ours. For example, [13] deals with quantile estimation based on the transnormal model.

High dimensional covariate issues have been important and intractable ones. However, some useful procedures have been proposed, for example, the SCAD in [9], the Lasso in [30], and the group Lasso in [36] and [26]. The properties of the Lasso were studied in [40] and [2]. The adaptive Lasso was proposed by [40] and it has the selection consistency property. The SCAD cannot deal with too many covariates and needs some screening procedures such as the SIS procedure in [11]. [15] proposed a quantile based screening procedure. There are some papers on screening procedures for varying coefficient and additive models, for example, [8], [10], and [20]. Forward type selection procedures are considered in e.g. [33] and [17]. We name [3], [14], and [32] as general references on high-dimensional issues.

Because parsimonious modelling is crucial for statistical analysis, simultaneous variable selection and structure identification in semiparametric regression models has been studied by many authors, see, among others, [37], [22], [35], [6], [23], and [16]. Another important reason to attain this purpose is that in some high-dimensional situations, there may be a lack of priori knowledge on how to decide which covariates to be included in the parametric part and which covariates to be included in the nonparametric part. On the other hand, to the best of our knowledge, no theoretical sound procedure has been proposed to achieve the aforementioned goal in the high-dimensional quantile regression setups. Note that [22] and [23] proposed using the estimated derivatives of coefficient functions to identify the structures of additive models. These estimated derivatives, however, usually have slow convergence rates. Moreover, as shown in Section S.3 of the

supplementary document, the conditions imposed on the B-spline basis functions in [22] and [23] seem too stringent to be satisfied in practice. Instead of relying on the estimated derivatives of coefficient functions, we appeal to the orthogonal decomposition method through introducing an orthonormal spline basis with desirable properties as in [16], which is devoted to the study of Cox regression models. Our approach not only can be justified theoretically under a set of reasonable assumptions, but also enables a unified analysis of varying coefficient models and additive models. The single index model is another important semiparametric quantile regression model. However, we don't deal with the model because the theoretical treatment is completely different from that of the varying coefficient and additive model. We just refer to [39] and [25] here.

The Lasso for quantile linear regression is considered in [1] and the adaptively weighted Lasso for quantile linear regression are considered in [7] and [38]. Some authors such as [18] and [29] deal with group Lasso procedures for additive models and varying coefficient models, respectively. [24] applied a reproducing kernel Hilbert space approach to additive models. [28] deals with SCAD type variable selection for parametric part. In [28], the authors applied the adaptively weighted Lasso iteratively to obtain their SCAD estimate starting from the Lasso estimate. However, in the quantile regression setup, there doesn't seem to exist any theoretical or numerical result for simultaneous variable selection and structure identification based on the adaptively weighted group Lasso, in particular when its penalty is determined by a data-driven fashion. To fill this gap, we establish selection consistency of AWG-Lasso and AWG-Lasso+HDIC in Section 3, and illustrate the finite sample performance of AWG-Lasso+HDIC through a simulation study in Section 4. Our simulation study reveals that AWG-Lasso+HDIC performs satisfactorily in terms of true positive and true negative rates.

This paper is organized as follows: We describe our procedures in Section 2. We present our theoretical results in Section 3. The results of numerical studies are given in Section 4. We state assumptions and prove our main results in Section 5 and describe some important properties of B-spline bases in the supplementary document, which also contains a real application of the proposed methods and more technical details.

We end this section with some notation used throughout the paper. \bar{A} and $|A|$ stand for the complement and the number of the elements of a set A , respectively. For a vector a , $|a|$ and a^T are the Euclidean norm and the transpose, respectively. For a function g on the unit interval, $\|g\|$ and $\|g\|_\infty$ stand for the L_2 and sup norms, respectively. We denote the maximum and minimum eigenvalues of a matrix A by $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$,

respectively. Besides, C, C_1, C_2, \dots , are generic positive constants and their values may change from line to line. Note that $a_n \sim b_n$ means $C_1 < a_n/b_n < C_2$ and that $a \vee b$ and $a \wedge b$ stand for the maximum and the minimum of a and b , respectively. Convergence in probability is denoted by \xrightarrow{p} .

2 Simultaneous variable selection and structure identification

We consider varying coefficient models and additive models in this paper. We can deal with both models in the same way and we concentrate on varying coefficient models in sections 2 and 3 to save space. We present the specific procedure for additive models in the supplement.

Suppose that we have n i.i.d. observations $\{(Y_i, \mathbf{X}_i, Z_i)\}_{i=1}^n$, where $\mathbf{X}_i = (X_{i1}, X_{i2}, \dots, X_{ip})^T$ is a p -dimensional covariate vector and Z_i is a scalar index covariate. Then we assume a quantile varying coefficient model holds for these observations. First we define the τ -th quantile check function $\rho_\tau(u)$ and its derivative $\rho'_\tau(u)$ by

$$\rho_\tau(u) = u(\tau - I\{u \leq 0\}) \quad \text{and} \quad \rho'_\tau(u) = \tau - I\{u \leq 0\}.$$

Then our varying coefficient model is

$$Y_i = \sum_{j=1}^p X_{ij} g_j(Z_i) + \epsilon_i, \tag{1}$$

where $Z_i \in [0, 1]$ and $E\{\rho'_\tau(\epsilon_i) | \mathbf{X}_i, Z_i\} = 0$. Usually we take $X_{i1} \equiv 1$ for varying coefficient models.

To deal with partially linear varying coefficient models, we decompose $g_j(z)$ as $g_j(z) = g_{cj} + g_{vj}(z)$, where

$$g_{cj} = \int_0^1 g_j(z) dz \quad \text{and} \quad g_{vj}(z) = g_j(z) - g_{cj}.$$

We define the index set, $\mathcal{S}^0 = (\mathcal{S}_c^0, \mathcal{S}_v^0)$, for the true model, where

$$\mathcal{S}_c^0 = \{j | g_{cj} \neq 0\} \quad \text{and} \quad \mathcal{S}_v^0 = \{j | g_{vj}(z) \not\equiv 0\}.$$

The index set for a candidate model can be similarly given by $\mathcal{S} = (\mathcal{S}_c, \mathcal{S}_v)$. In the following, we refer to \mathcal{S}^0 and \mathcal{S} as the true model and the candidate model, respectively whenever confusion is unlikely. When some j 's satisfy both $j \in \mathcal{S}_c^0$ and

$j \notin \mathcal{S}_v^0$ simultaneously, our true model is a partially linear varying coefficient model, for example, $\mathcal{S}^0 = (\{1, 2, 3\}, \{1, 2\})$ with $\mathcal{S}_c^0 = \{1, 2, 3\}$ and $\mathcal{S}_v^0 = \{1, 2\}$. Moreover, $\mathcal{S}_1 \supset \mathcal{S}_2$ means $\mathcal{S}_{c1} \supset \mathcal{S}_{c2}$ and $\mathcal{S}_{v1} \supset \mathcal{S}_{v2}$, where $\mathcal{S}_j = (\mathcal{S}_{cj}, \mathcal{S}_{vj})$, $j = 1, 2$. In addition, $\mathcal{S}_1 \cup \mathcal{S}_2 = (\mathcal{S}_{c1} \cup \mathcal{S}_{c2}, \mathcal{S}_{v1} \cup \mathcal{S}_{v2})$.

We use the regression spline method to estimate coefficient functions and the covariates for regression spline are defined by

$$\mathbf{W}_i = \mathbf{X}_i \otimes \mathbf{B}(Z_i), \quad (2)$$

where $\mathbf{B}(z) = (B_1(z), B_2(z), \dots, B_L(z))^T$ is an orthonormal basis constructed from the equispaced B-spline basis $\mathbf{B}_0(z) = (B_{01}(z), \dots, B_{0L}(z))^T$ on $[0, 1]$ and \otimes is the Kronecker product. We can represent $\mathbf{B}(z)$ as $\mathbf{B}(z) = A_0 \mathbf{B}_0(z)$ and we calculate the $L \times L$ matrix A_0 numerically. As in [16], let $\mathbf{B}(z)$ satisfy $B_1(z) = 1/\sqrt{L}$, $B_2(z) = \sqrt{12/L}(z - 1/2)$, and

$$\int_0^1 \mathbf{B}(z)(\mathbf{B}(z))^T dz = L^{-1} I_L. \quad (3)$$

We denote the $L \times L$ identity matrix by I_L . Note that $B_1(z)$ is for g_{cj} (the j -th constant component) and $\mathbf{B}_{-1}(z) = (B_2(z), \dots, B_L(z))^T$ is for $g_{vj}(z)$ (the j -th non-constant component). More details are given in Section S.3 of the supplement.

To carry out simultaneous variable selection and structure identification, we apply AWG-Lasso to

$$Y_i = \mathbf{W}_i^T \boldsymbol{\gamma} + \epsilon'_i, \quad (4)$$

where $\boldsymbol{\gamma} = (\boldsymbol{\gamma}_1^T, \dots, \boldsymbol{\gamma}_p^T)^T$. For a given $\lambda > 0$, the corresponding objective function is given by

$$Q_V(\boldsymbol{\gamma}; \lambda) = \frac{1}{n} \sum_{i=1}^n \rho_\tau(Y_i - \mathbf{W}_i^T \boldsymbol{\gamma}) + \lambda \sum_{j=1}^p (w_{1j} |\gamma_{1j}| + w_{-1j} |\gamma_{-1j}|), \quad (5)$$

where $\{(w_{1j}, w_{-1j})\}_{j=1}^p$ is obtained from some initial estimates such as Lasso and group Lasso, and $(\gamma_{1j}, \gamma_{-1j})^T = \boldsymbol{\gamma}_j$, noting that γ_{1j} is for $B_1(z)$ and γ_{-1j} is for $\mathbf{B}_{-1}(z)$. Minimizing $Q_V(\boldsymbol{\gamma}; \lambda)$ w.r.t. $\boldsymbol{\gamma}$, one gets

$$\hat{\boldsymbol{\gamma}}^\lambda = \underset{\boldsymbol{\gamma} \in \mathbb{R}^{pL}}{\operatorname{argmin}} Q_V(\boldsymbol{\gamma}; \lambda).$$

Denote $\hat{\boldsymbol{\gamma}}^\lambda$ by $(\hat{\gamma}_{11}^\lambda, \hat{\gamma}_{-11}^{\lambda T}, \dots, \hat{\gamma}_{1p}^\lambda, \hat{\gamma}_{-1p}^{\lambda T})^T$. Then, the model selected by AWG-Lasso is $\hat{\mathcal{S}}^\lambda = (\hat{\mathcal{S}}_c^\lambda, \hat{\mathcal{S}}_v^\lambda)$, where $\hat{\mathcal{S}}_c^\lambda = \{j \mid \hat{\gamma}_{1j}^\lambda \neq 0\}$ and $\hat{\mathcal{S}}_v^\lambda = \{j \mid \hat{\gamma}_{-1j}^\lambda \neq \mathbf{0}\}$, and this enables us to identify variables and structures simultaneously.

Theorem 1 in Section 3 establishes the selection consistency of $\widehat{\mathcal{S}}^\lambda$ under a set of general conditions on $\{(w_{1j}, w_{-1j})\}_{j=1}^p$ and a strong sparsity condition on the regression coefficients that $|\mathcal{S}_c^0|$ and $|\mathcal{S}_v^0|$ are bounded. Theorem 1, however, also requires that λ falls into a suitable interval, which can sometimes be hard to decide in practice. We therefore introduce a BIC-type criterion, HDIC, to choose a λ in a data-driven fashion. Express \mathbf{W}_i as $(v_{11i}, \mathbf{v}_{-11i}^T, \dots, v_{1pi}, \mathbf{v}_{-1pi}^T)^T$, where $(v_{1ji}, \mathbf{v}_{-1ji}^T)^T$ is the regressor vector corresponding to γ_j . For a given model $\mathcal{S} = (\mathcal{S}_c, \mathcal{S}_v)$, define $R_V(\boldsymbol{\gamma}_\mathcal{S})$ and $\tilde{\boldsymbol{\gamma}}_\mathcal{S}$ by

$$R_V(\boldsymbol{\gamma}_\mathcal{S}) = \frac{1}{n} \sum_{i=1}^n \rho_\tau(Y_i - \mathbf{W}_{i\mathcal{S}}^T \boldsymbol{\gamma}_\mathcal{S}) \quad \text{and} \quad \tilde{\boldsymbol{\gamma}}_\mathcal{S} = \underset{\boldsymbol{\gamma}_\mathcal{S} \in R^{|\mathcal{S}_c|+(L-1)|\mathcal{S}_v|}}{\operatorname{argmin}} R_V(\boldsymbol{\gamma}_\mathcal{S}), \quad (6)$$

where $\mathbf{W}_{i\mathcal{S}} \in R^{|\mathcal{S}_c|+(L-1)|\mathcal{S}_v|}$ consists of $\{v_{1ji} \mid j \in \mathcal{S}_c\}$ and $\{\mathbf{v}_{-1ji} \mid j \in \mathcal{S}_v\}$. The corresponding coefficient vector $\boldsymbol{\gamma}_\mathcal{S}$ consists of $\{\gamma_{1ji} \mid j \in \mathcal{S}_c\}$ and $\{\gamma_{-1ji} \mid j \in \mathcal{S}_v\}$ as well. The elements of these vectors are suitably arranged. In this paper, we sometimes take two index sets \mathcal{S}_1 and \mathcal{S}_2 satisfying $\mathcal{S}_1 \subset \mathcal{S}_2$ and compare $\boldsymbol{\gamma}_{\mathcal{S}_1}$ and $\boldsymbol{\gamma}_{\mathcal{S}_2}$ by enlarging $\boldsymbol{\gamma}_{\mathcal{S}_1}$ with 0 elements or something, for example, $(\boldsymbol{\gamma}_{\mathcal{S}_1}^T, \mathbf{0}^T)^T$. Then $(\boldsymbol{\gamma}_{\mathcal{S}_1}^T, \mathbf{0}^T)^T$ and $\boldsymbol{\gamma}_{\mathcal{S}_2}$ have the same dimension and the elements of these vectors are assumed to be conformably rearranged.

The HDIC value for model \mathcal{S} is stipulated by

$$\text{HDIC}(\mathcal{S}) = \log R_V(\tilde{\boldsymbol{\gamma}}_\mathcal{S}) + (|\mathcal{S}_c| + (L-1)|\mathcal{S}_v|) \frac{q_n \log p_n}{2n}, \quad (7)$$

where $p_n = p \vee n$ and $q_n \rightarrow \infty$ at a slow rate described in Section 5. We consider a set of models $\{\widehat{\mathcal{S}}^\lambda\}$ chosen by AWG-Lasso, where $\lambda \in \Lambda$ with Λ being a prescribed set of positive numbers, and select $\widehat{\mathcal{S}}^{\hat{\lambda}}$ among $\{\widehat{\mathcal{S}}^\lambda\}$, where

$$\hat{\lambda} = \underset{\lambda \in \Lambda, |\widehat{\mathcal{S}}_c^\lambda| \leq M_c, |\widehat{\mathcal{S}}_v^\lambda| \leq M_v}{\operatorname{argmin}} \text{HDIC}(\widehat{\mathcal{S}}^\lambda),$$

with M_c and M_v being known upper bounds for $|\mathcal{S}_c^0|$ and $|\mathcal{S}_v^0|$, respectively. Under some regularity conditions, the consistency of $\widehat{\mathcal{S}}^{\hat{\lambda}}$ is established in Corollary 1.

Note that in the case of high-dimensional sparse linear models, it is shown in [17] that (7) with $\rho_\tau(\cdot)$ replaced by the squared loss $(\cdot)^2$ can be used in conjunction with the orthogonal greedy algorithm (OGA) to yield selection consistency. The major difference between (7) and the BIC-type criteria considered in [21] is that we deal with semiparametric models in this paper. It seems difficult to derive the consistency of $\widehat{\mathcal{S}}^{\hat{\lambda}}$ in any high-dimensional regression setups without the additional penalty term q_n in (7).

3 Consistency results

We prove the consistency of AWG-Lasso and AWG-Lasso+HDIC separately in Subsection 3.1 and 3.2. It is worth pointing out that due to the similarity between (4)-(7) and (S.2)-(S.5) in the supplement, the theoretical treatment is almost the same for the two types of models considered in this paper. Therefore, this section concentrates only on the varying coefficient model. On the other hand, our numerical studies are conducted for both types of models, see Section 4.

3.1 Adaptively weighted group Lasso

The consistency of AWG-Lasso for suitably chosen λ and weights is stated in Theorem 1. The proof of Theorem 1 is reliant on the methods of [7], [38], and [28] subject to non-trivial modifications. The details are deferred to Section 5. For clarity of presentation, all the technical assumptions of Theorem 1 are also given in Section 5. Roughly speaking, we assume that the coefficient functions have second order derivatives and we put $L = c_L n^{1/5}$. More smoothness is necessary for Theorem 2. If X_{ij} is uniformly bounded, the Hölder continuity of the second order derivatives with exponent $\alpha = 1/2$ is sufficient for Theorem 2.

Define $d_V(\mathcal{S}) = |\mathcal{S}_c| + (L - 1)|\mathcal{S}_v|$ and let $w_{\mathcal{S}^0}$ denote a weight vector consisting of $\{w_{1j} \mid j \in \mathcal{S}_c^0\}$ and $\{w_{-1j} \mid j \in \mathcal{S}_v^0\}$. For an index set \mathcal{S} , we define $\widehat{\gamma}_{\mathcal{S}}^\lambda$ by

$$\widehat{\gamma}_{\mathcal{S}}^\lambda = \operatorname{argmin}_{\gamma_{\mathcal{S}} \in R^{d_V(\mathcal{S})}} Q_V(\gamma_{\mathcal{S}}; \lambda).$$

Then $\widehat{\gamma}_{\mathcal{S}^0}^\lambda$ is an oracle estimator on $R^{d_V(\mathcal{S}^0)}$ with the knowledge of \mathcal{S}^0 . Assumption A2 assumes that the relevant coefficients and the coefficient functions are large enough to be detected.

Theorem 1 *Assume that Assumptions A1, A3-5 and B1-4 in Section 5 hold. Moreover, assume*

$$\max_{j \in \mathcal{S}_c^0} w_{1j} \vee \max_{j \in \mathcal{S}_v^0} w_{-1j} = O_p(1), \tag{8}$$

and for some sufficiently large $0 < a_1, a_2 < \infty$,

$$\min_{j \notin \mathcal{S}_c^0} w_{1j} \geq (a_1 |w_{\mathcal{S}^0}|) \vee 1 \quad \text{and} \quad \min_{j \notin \mathcal{S}_v^0} w_{-1j} \geq (a_2 |w_{\mathcal{S}^0}|) \vee 1, \tag{9}$$

with probability tending to 1. We enlarge $\widehat{\gamma}_{\mathcal{S}^0}$ by adding 0 elements for the \mathcal{S}^{0c} part so that $(\widehat{\gamma}_{\mathcal{S}^0}^{\lambda T}, \mathbf{0}^T)^T \in R^{pL}$ and define $\widehat{\mathcal{S}}^\lambda$ from this $(\widehat{\gamma}_{\mathcal{S}^0}^{\lambda T}, \mathbf{0}^T)^T$. Then for any λ satisfying

$$a_3 \frac{(\log p_n)^{1/2}}{n^{1/2}} \leq \lambda \leq (\log n)^\kappa \frac{(\log p_n)^{1/2}}{n^{1/2}} \quad (10)$$

asymptotically, where a_3 is a sufficiently large constant and κ is any positive constant, $(\widehat{\gamma}_{\mathcal{S}^0}^{\lambda T}, \mathbf{0}^T)^T (= \widehat{\mathcal{S}}^\lambda)$ is actually an optimal solution to minimizing $Q_V(\gamma; \lambda)$ w.r.t. $\gamma \in R^{pL}$ with probability tending to 1. If Assumption A2 also holds, we have for $\widehat{\mathcal{S}}^\lambda$ defined here that

$$\lim_{n \rightarrow \infty} \mathbb{P}(\widehat{\mathcal{S}}^\lambda = \mathcal{S}^0) = 1.$$

The order of $L^{1/2}\lambda$ from (10) is the standard one in the literature since $(\log p_n)^{1/2}$ is due to the large number of covariates and $(L/n)^{1/2}$ is the standard rate for regression spline estimation. Recall that our normalization factor of the orthonormal basis is $1/L$. The upper bound of λ in Theorem 1 is a technical one since we approximate $R_V(\gamma)$ by a quadratic function in γ on a suitable bounded region.

We will further discuss the convergence rate of the AWG-Lasso estimators and present two examples of data-driven weights.

First we discuss the convergence rate of the AWG-Lasso estimators by referring to Proposition 1 in Section 5. We have derived the consistency of $\widehat{\mathcal{S}}^\lambda$ in Theorem 1. Then if we apply Proposition 1 with $\mathcal{S} = \mathcal{S}^0$, we have from Remark 1 there that

$$\mathbb{P}(|\widehat{\gamma}_{\mathcal{S}}^\lambda - \gamma_{\mathcal{S}}^*| \geq \eta_n) \rightarrow 0,$$

where $\eta_n \sim L\{(n^{-1} \log p_n)^{1/2} + \lambda|w_{\mathcal{S}^0}|\}$. We state the proposition for the proofs of Theorems 1 and 2 to take care of uniformity with respect to the indices of covariates and we can improve the rate slightly and replace $\log p_n$ with $\log n$ for this one index set \mathcal{S}^0 . Hence the convergence rate of the oracle AWG-Lasso estimators of g_{cj} , $j \in \mathcal{S}_c^0$, and g_{vj} , $j \in \mathcal{S}_v^0$, is $L^{1/2}\{(n^{-1} \log n)^{1/2} + \lambda|w_{\mathcal{S}^0}^0|\}$ in the setup of Remark 1.

Next we present two examples of data-driven weights here. A simple sufficient condition for (9) is that with probability tending to 1,

$$\frac{\min_{j \notin \mathcal{S}_c^0} w_{1j} \wedge \min_{j \notin \mathcal{S}_v^0} w_{-1j}}{1 \vee \max_{j \in \mathcal{S}_c^0} w_{1j} \vee \max_{j \in \mathcal{S}_v^0} w_{-1j}} \rightarrow \infty. \quad (11)$$

Example 1(Adaptive Lasso type weights). We need an initial estimator denoted by $\bar{\gamma} = (\bar{\gamma}_{11}, \bar{\gamma}_{-11}^T, \dots, \bar{\gamma}_{1p}, \bar{\gamma}_{-1p}^T)^T$ from the group Lasso as in [29] and [18]. Note that

$L^{-1/2}|\bar{\gamma}_{1j}|$ and $L^{-1/2}|\bar{\gamma}_{-1j}|$ from [29] and [18] are consistent estimates of $|g_{cj}|$ and $\|g_{vj}\|$, respectively. Actually they have the convergence rates smaller than $CL^{1/2}\lambda$ for some sufficiently large C and λ in Theorem 1. Hence

$$w_{1j} = (L^{-1/2}|\bar{\gamma}_{1j}|)^{-\eta} \quad \text{and} \quad w_{-1j} = (L^{-1/2}|\bar{\gamma}_{-1j}|)^{-\eta} \quad (12)$$

satisfy the conditions (8) and (9) for any positive fixed η if we have for some positive C that $\min_{j \in \mathcal{S}_c^0} |g_{cj}| \wedge \min_{j \in \mathcal{S}_v^0} \|g_{vj}\| > C$. On the other hand, if $\min_{j \in \mathcal{S}_c^0} |g_{cj}| \wedge \min_{j \in \mathcal{S}_v^0} \|g_{vj}\| \rightarrow 0$ slowly as in Assumption A2 in Section 5, we can cope with this situation theoretically by making a suitable adjustment to the order of λ . Note that $\lambda w_{1j} = (\xi_n \lambda)(\xi_n^{-1} w_{1j})$ and $\lambda w_{-1j} = (\xi_n \lambda)(\xi_n^{-1} w_{-1j})$ for a suitable ξ_n and that $\xi_n \lambda$, $\xi_n^{-1} w_{1j}$, and $\xi_n^{-1} w_{-1j}$ have only to meet the assumptions in Theorem 1. However, we usually have no knowledge of the order of $\min_{j \in \mathcal{S}_c^0} |g_{cj}| \wedge \min_{j \in \mathcal{S}_v^0} \|g_{vj}\|$ in advance and this kind of adjustment to λ may be practically difficult. Or then we should try a very wide range of λ .

Example 2 (SCAD-based weights). With the initial estimator $\bar{\gamma}$ obtained from the Lasso penalty estimators such as in [29] and [18], we apply one-step LLA (local linear approximation) to the SCAD penalty as in [12] to obtain $\{(w_{1j}, w_{-1j})\}$. More specifically, we set

$$\lambda w_{1j} |\gamma_{1j}| = p'_{\lambda L^{1/2}}(L^{-1/2}|\bar{\gamma}_{1j}|)(L^{-1/2}|\gamma_{1j}|) \quad \text{and} \quad (13)$$

$$\lambda w_{-1j} |\gamma_{-1j}| = p'_{\lambda L^{1/2}}(L^{-1/2}|\bar{\gamma}_{-1j}|)(L^{-1/2}|\gamma_{-1j}|), \quad (14)$$

where $p_\lambda(\cdot)$ is the SCAD penalty function. Some authors as [28] applied this kind of AGW-Lasso iteratively to calculate their SCAD estimates.

Because of the properties of the SCAD penalty function, there are positive constants C_1 , C_2 , and C_3 such that if with probability tending to 1,

$$\frac{\min_{j \in \mathcal{S}_c^0} L^{-1/2}|\bar{\gamma}_{1j}| \wedge \min_{j \in \mathcal{S}_v^0} L^{-1/2}|\bar{\gamma}_{-1j}|}{\lambda L^{1/2}} > C_1 \quad \text{and} \quad (15)$$

$$\frac{\max_{j \notin \mathcal{S}_c^0} L^{-1/2}|\bar{\gamma}_{1j}| \vee \max_{j \notin \mathcal{S}_v^0} L^{-1/2}|\bar{\gamma}_{-1j}|}{\lambda L^{1/2}} < C_2, \quad (16)$$

then we have with probability tending to 1,

$$w_{1j} = 0(j \in \mathcal{S}_c^0) \quad \text{and} \quad w_{-1j} = 0(j \in \mathcal{S}_v^0) \quad \text{and} \quad w_{1j} > C_3(j \notin \mathcal{S}_c^0) \quad \text{and} \quad w_{-1j} > C_3(j \notin \mathcal{S}_v^0).$$

Thus the weights given in (13) and (14) obey (8) and (9). If necessary, we multiply λ and the weights by $1/C_4$ and C_4 , respectively, where C_4 is a sufficiently large constant

and this adjustment does not essentially affect the condition (10). If

$$\frac{\min_{j \in \mathcal{S}_c^0} |g_{cj}| \wedge \min_{j \in \mathcal{S}_v^0} \|g_{vj}\|}{\lambda L^{1/2}} \rightarrow \infty,$$

we will have (15) and (16). Note that these weights don't meet (11).

3.2 Consistency of AWG-Lasso+HDIC

To state the main result of this subsection, we need to introduce Assumption A1, which assumes that $|\mathcal{S}_c^0| \leq C_c$ and $|\mathcal{S}_v^0| \leq C_v$ for some fixed C_c and C_v . Let M_c and M_v be known positive integers fixed with n such that $C_c < M_c$ and $C_v < M_v$. Define

$$\widehat{\mathcal{S}} = \underset{|\mathcal{S}_c| \leq M_c \text{ and } |\mathcal{S}_v| \leq M_v}{\operatorname{argmin}} \operatorname{HDIC}(\mathcal{S}).$$

Under certain regularity conditions, the next theorem and corollary show that both $\widehat{\mathcal{S}}$ and $\widehat{\mathcal{S}}^\lambda$ are consistent estimates of \mathcal{S}^0 . We need to replace Assumptions A2–5 and B1–4 with Assumptions A2'–A5' and B1'–B4' to carry out subtle evaluations of $R_V(\gamma_{\mathcal{S}})$ in the proof since we deal with high-dimensional semiparametric models. All the technical assumptions of Theorem 2 are also given in Section 5.

Theorem 2 *Assume that Assumptions A1, A2'–A5', B1'–B4' and B5 in Section 5 hold. Then,*

$$\lim_{n \rightarrow \infty} P(\widehat{\mathcal{S}} = \mathcal{S}^0) = 1.$$

Corollary 1 *We assume the same assumptions as in Theorem 2 and that (8) and (9) hold true. Then for Λ satisfying $\Lambda \subset [c_n^{-1} \sqrt{\log p_n/n}, c_n \sqrt{\log p_n/n}]$ and $\{c_n \sqrt{\log p_n/n}\} \in \Lambda$, where $c_n \rightarrow \infty$ and $c_n/(\log n)^\kappa \rightarrow 0$ for some $\kappa > 0$, we have*

$$\lim_{n \rightarrow \infty} P(\widehat{\mathcal{S}}^\lambda = \mathcal{S}^0) = 1.$$

Some comments are in order. While $\widehat{\mathcal{S}}$ can achieve selection consistency without the help of AWG-Lasso, it seems difficult to obtain $\widehat{\mathcal{S}}$ directly when p is large and M_c and M_v are not very small. On the other hand, $\widehat{\mathcal{S}}^\lambda$ is applicable in most practical situations. We also note that Theorem 2 extends the result in [21] and can be viewed as a generalization of the BIC result in [38] to the semiparametric setup, which is of fundamental interest from both theoretical and practical perspectives. Like [38], [19] also confines its attention to linear quantile models. Moreover, it seems difficult to

extend the proof in [19] to situations where the dimension of the true model tends to infinity. Finally, we mention that there is another version of HDIC,

$$\text{HDIC}_{\text{II}}(\mathcal{S}) = R_V(\tilde{\gamma}_{\mathcal{S}}) + (|\mathcal{S}_c| + (L - 1)|\mathcal{S}_v|) \frac{q_n \log p_n}{2n}, \quad (17)$$

which becomes

$$\text{HDIC}_{\text{II}}(\mathcal{S}) = R_V(\tilde{\gamma}_{\mathcal{S}}) + (|\mathcal{S}_l| + (L - 2)|\mathcal{S}_a|) \frac{q_n \log p_n}{2n} \quad (18)$$

in the case of additive models. It can be shown that HDIC_{II} and HDIC share the same asymptotic properties and their finite sample performance will be compared in the next section.

4 Numerical studies

In this section, we evaluate the performance of AWG-Lasso+HDIC and AWG-Lasso+ HDIC_{II} using one varying coefficient model and two additive models in the case of $pL > n$. We set $q_n = 1$ in these numerical studies since the optimal choice of q_n in finite sample remains unsettled and is worth further investigation. Moreover, $\{(w_{1j}, w_{-1j})\}$ in (5) are assigned according to (13) and (14), and $\{(w_{2j}, w_{-2j})\}$ in (S.3) are determined in a similar fashion.

In our simulation study, we consider one varying coefficient model (Example 1) and two additive models (Examples 2 and 3). In these examples, we set $(n, p) = (500, 400)$, $L = 6$, $\tau = 0.5$, $M_c = M_v = M_l = M_a = 20$ and

$$\Lambda = \{c_n^{-1} \sqrt{\log p/n} + kd_n, k = 1, \dots, 50\},$$

where $c_n = 2 \log n$ and $d_n = \{(c_n - c_n^{-1}) \sqrt{\log p/n}\}/50$.

Based on a $\lambda \in \Lambda$ and the weights described above, we employ the alternating direction method of multipliers (ADMM) to minimize (5) ((S.3)) over γ (γ_{-1}), and then choose the λ minimizing $\text{HDIC}(\hat{\mathcal{S}}^\lambda)$ defined in (7) ((S.5)) over $\lambda \in \Lambda$, and the λ minimizing $\text{HDIC}_{\text{II}}(\hat{\mathcal{S}}^\lambda)$ defined in (17) ((18)) over the same set. We conduct 50 simulations and the performance of AWG-Lasso+HDIC and AWG-Lasso+ HDIC_{II} in Examples 1–3 is documented in Tables 1–3, respectively. For the purpose of comparison, we also use the Rqpen package in R (see cv.rq.group.pen) to implement the group Lasso method in Example 1–3. In addition, the adaptive group Lasso method introduced in [29] for varying coefficient models (referred to as the T-method), and the group Lasso

method introduced in [18] for additive models (referred to as the K-method) are included. Note that since our goal is to identify structures in addition to selecting variables, these three methods are conducted based on the orthonormal basis functions proposed in this paper, which enable one to distinguish between constant and non-constant components for varying coefficient models (or liner and non-linear components for additive models). On the other hand, we use their original penalties, not the divided ones like ours. The performance of these three methods is also presented in Tables 1–3. In the the Rqpen package, the L_1 norm is used instead of the L_2 norm inside the penalty functions. See the document for the details. This may be the cause of different performances from the other methods.

Example 1. We generate the output variables Y_1, \dots, Y_n using the varying coefficient model,

$$Y_i = \sum_{j=1}^p X_{ij} g_j(Z_i) + \epsilon_i,$$

where ϵ_i , Z_i and $\{X_{ij}\}_{j=1}^p$ are independently generated from $N(0, 0.5^2)$, $U(0, 1)$ and $U(0, 100)$ distributions, respectively. Following [16], the coefficient functions $g_j(z)$ are set to

$$g_1(z) = g_2(z) = 1, g_3(z) = 4z, g_4(z) = 4z^2, g_j(z) = 0, \quad 5 \leq j \leq p.$$

Therefore, X_{i1} and X_{i2} are relevant covariates with constant coefficients, X_{i3} and X_{i4} are relevant covariates with non-constant coefficients, whereas X_{i5}, \dots, X_{ip} , are irrelevant variables. Since our goal is to identify both relevant variables and the structures of relevant coefficients, define

$$C_{sj} = I_{\{g_j(\cdot) \text{ is identified as a constant function at the } s\text{th replication}\}},$$

$$NC_{sj} = I_{\{g_j(\cdot) \text{ is identified as a non-constant function at the } s\text{th replication}\}},$$

$$NS_{sj} = I_{\{g_j(\cdot) \text{ is identified as a zero function at the } s\text{th replication}\}}.$$

It is clear that $C_{sj} + NC_{sj} + NS_{sj} = 1$ for each $1 \leq j \leq p$. We further define the true negative rate (TNR) and the strictly true positive rate (STPR),

$$\text{TNR}_s = \frac{\sum_{j=5}^p I_{\{NS_{sj}=1\}}}{p-4} \quad \text{and} \quad \text{STPR}_s = \frac{\sum_{j=1}^2 I_{\{C_{sj}=1\}} + \sum_{j=3}^4 I_{\{NC_{sj}=1\}}}{4},$$

noting that $\text{STPR}_s = 1$ if at the s th replication, X_{i1} and X_{i2} are identified as relevant variables with constant coefficients and X_{i3} and X_{i4} are identified as relevant variables

with non-constant coefficients. Therefore, STPR_s can be viewed as a stringent version of the conventional true positive rate, which treats constant and non-constant coefficient functions indifferently. Now, the performance measures of a selection method are specified as follows:

$$\begin{aligned} C_j &= \frac{1}{50} \sum_{s=1}^{50} C_{sj}, \quad \text{NC}_j = \frac{1}{50} \sum_{s=1}^{50} \text{NC}_{sj}, \quad \text{NS}_j = \frac{1}{50} \sum_{s=1}^{50} \text{NS}_{sj}, \\ \text{TNR} &= \frac{1}{50} \sum_{s=1}^{50} \text{TNR}_s, \quad \text{STPR} = \frac{1}{50} \sum_{s=1}^{50} \text{STPR}_s. \end{aligned}$$

The performance of AWG-Lasso+HDIC, AWG-Lasso+HDIC_{II}, Rqpen, and T-method on $(C_j, \text{NC}_j, \text{NS}_j), j = 1, \dots, 4$, STPR and TNR is demonstrated in Table 1. Table 1 shows that AWG-Lasso+HDIC and AWG-Lasso+HDIC_{II} have high capability in identifying the true variables and true structures in the sense that $C_1=C_2=\text{NC}_3=\text{NC}_4=\text{STPR}=1$ hold for the two methods. Table 1 also reveals that both methods perform satisfactorily in identifying irrelevant variables since their TNR values are quite close to 1. Because Rqpen encounters singularity problems in many replications, its performance measures are set to missing in Table 1. The T-method performs quite well in identifying irrelevant variables and non-constant functions because its TNR, NC_3 , and NC_4 are equal to 1. The method, however, erroneously treats constant functions as non-constant ones, leading to a low STPR value of 0.5.

Example 2. We generate Y_1, \dots, Y_n from the following additive model,

$$Y_i = \mu + \sum_{j=1}^p g_j(X_{ij}) + \epsilon_i, \quad (19)$$

where $\mu = 0$, ϵ_i and $\{X_{ij}\}_{j=1}^p$ follow $N(0, 0.5^2)$ and $U(0, 1)$, respectively. Following [16] again, we set

$$\begin{aligned} g_1(x) &= g_2(x) = 2^{1/2}(x - 1/2), \quad g_3(x) = 2^{-1/2} \cos(2\pi x) + (x - 1/2), \\ g_4(x) &= \sin(2\pi x), \quad g_i(x) = 0, \quad 5 \leq i \leq p, \end{aligned} \quad (20)$$

noting that X_{i1} and X_{i2} are relevant through the linear functions $g_1(\cdot)$ and $g_2(\cdot)$, whereas X_{i3} and X_{i4} are relevant through the nonlinear functions $g_3(\cdot)$ and $g_4(\cdot)$. Let NS_{sj} and TNR_s be defined as in Example 1, and define

$$\begin{aligned} L_{sj} &= I_{\{g_j(\cdot) \text{ is identified as a linear function at the } s\text{th replication}\}}, \\ \text{NL}_{sj} &= I_{\{g_j(\cdot) \text{ is identified as a non-linear function at the } s\text{th replication}\}}, \\ \text{STPR}_s &= \frac{\sum_{j=1}^2 I_{\{L_{sj}=1\}} + \sum_{j=3}^4 I_{\{\text{NL}_{sj}=1\}}}{4}. \end{aligned}$$

Then, the performance measures of AWG-Lasso+HDIC, AWG-Lasso+HDIC_{II}, Rqpen, and K-method are given by

$$\begin{aligned} L_j &= \frac{1}{50} \sum_{s=1}^{50} L_{sj}, \quad NL_j = \frac{1}{50} \sum_{s=1}^{50} NL_{sj}, \quad NS_j = \frac{1}{50} \sum_{s=1}^{50} NS_{sj}, \\ TNR &= \frac{1}{50} \sum_{s=1}^{50} TNR_s, \quad STPR = \frac{1}{50} \sum_{s=1}^{50} STPR_s, \end{aligned}$$

and summarized in Tables 2. Table 2 shows that $L_1 = L_2 = 1$ hold for AWG-Lasso+HDIC, AWG-Lasso+HDIC_{II}, and Rqpen, implying that these three methods can easily identify relevant linear functions. In addition, the NL_3 and NL_4 of these three methods are equal (or close) to 1, leading to very high STRP values. While the TNR values of AWG-Lasso+HDIC and AWG-Lasso+HDIC_{II} are still very close to 1, Rqpen has a low TNR value of 0.67, revealing that the method may suffer from overfitting. On the other hand, the K-method can avoid overfitting and has the highest possible TNR value of 1. Moreover, its NL_3 and NL_4 are equal to 1, showing a good ability to identify non-linear functions. Unfortunately, the method fails to identify linear functions, resulting a low STPR value of 0.5.

Example 3. Suppose that Y_1, \dots, Y_n are still generated from model (19), but with (20) replaced by

$$\begin{aligned} g_1(x) &= \frac{3 \sin(2\pi x)}{(2 - \sin(2\pi x))} - 0.4641016, \quad g_2(x) = 6x(1 - x) - 1, \quad g_3(x) = 2x - 1, \\ g_4(x) &= x - 0.5, \quad g_5(x) = -x + 0.5, \quad g_i(x) = 0, \quad 6 \leq i \leq p, \end{aligned} \quad (21)$$

which are suggested in [22]. As observed in (21), X_{i1} and X_{i2} are relevant through the nonlinear functions $g_1(\cdot)$ and $g_2(\cdot)$, and $X_{i3} \sim X_{i5}$ are relevant through the linear functions $g_3(\cdot) \sim g_5(\cdot)$. With

$$TNR_s = \frac{\sum_{j=6}^p I_{\{NS_{sj}=1\}}}{p - 5} \quad \text{and} \quad STPR_s = \frac{\sum_{j=1}^2 I_{\{NL_{sj}=1\}} + \sum_{j=3}^5 I_{\{L_{sj}=1\}}}{5},$$

the performance measures of the methods considered in Example 2 are given by

$$\begin{aligned} L_j &= \frac{1}{50} \sum_{s=1}^{50} L_{sj}, \quad NL_j = \frac{1}{50} \sum_{s=1}^{50} NL_{sj}, \quad NS_j = \frac{1}{50} \sum_{s=1}^{50} NS_{sj}, \\ TNR &= \frac{1}{50} \sum_{s=1}^{50} TNR_s, \quad STPR = \frac{1}{50} \sum_{s=1}^{50} STPR_s, \end{aligned}$$

and summarized in Table 3. Table 3 shows that $NL_1 = NL_2 = L_3 = L_4 = L_5 = STPR = 1$ hold for AWG-Lasso+HDIC and AWG-Lasso+HDIC_{II}, suggesting that the two methods can perfectly identify the relevant variables as well as the corresponding functional structures. The two methods are also good at identifying irrelevant variables in terms of TNR values. The performance of the K-method in this example resembles that in Example 2. Rqpen still encounters overfitting as in Example 2. Moreover, it has a limited ability to identify linear functions although it can perfectly identify non-linear ones.

In conclusion, we note that the results of this section, together with those obtained in the previous sections, demonstrate that AWG-Lasso+HDIC and AWG-Lasso+HDIC_{II} have a strong ability to simultaneously identify the relevant (or irrelevant) variables and their corresponding structures in the high-dimensional quantile regression setup, a feature rarely reported in the literature. While the T- and K-methods also perform well in identifying relevant (or irrelevant) variables, they are not very successful in structure identification. This is mainly because the two methods don't penalize constant/linear and non-constant/non-linear terms separately. Rqpen can encounter numerical difficulties in high-dimensional varying coefficient models as demonstrated in Example 1. The performance of Rqpen in structure identification is as good as our method in Example 2, and slightly better than the K-method in Example 3. The method, however, often suffers from overfitting.

Table 1: $(C_i, NC_i, NS_i), i = 1, \dots, 4$, STPR, and TNR in Example 1

	$(n, p) = (500, 400)$				STPR	TNR
	(C_1, NC_1, NS_1)	(C_2, NC_2, NS_2)	(C_3, NC_3, NS_3)	(C_4, NC_4, NS_4)		
AWG-Lasso+HDIC	(1.0, 0.0, 0.0)	(1.0, 0.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	1.0	0.963
AWG-Lasso+HDIC _{II}	(1.0, 0.0, 0.0)	(1.0, 0.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	1.0	0.963
Rqpen	(-, -, -)	(-, -, -)	(-, -, -)	(-, -, -)	-	-
T-method	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	0.5	1.0

Table 2: $(L_i, NL_i, NS_i), i = 1, \dots, 4$, STPR, and TNR in Example 2

	$(n, p) = (500, 400)$				STPR	TNR
	(L_1, NL_1, NS_1)	(L_2, NL_2, NS_2)	(L_3, NL_3, NS_3)	(L_4, NL_4, NS_4)		
AWG-Lasso+HDIC	(1.0, 0.0, 0.0)	(1.0, 0.0, 0.0)	(0.0, 0.96, 0.04)	(0.0, 1.0, 0.0)	1.0	0.997
AWG-Lasso+HDIC _{II}	(1.0, 0.0, 0.0)	(1.0, 0.0, 0.0)	(0.0, 0.98, 0.02)	(0.02, 0.98, 0.0)	0.99	0.998
Rqpen	(1.0, 0.0, 0.0)	(1.0, 0.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	1.0	0.674
K-method	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	0.5	1.0

Table 3: $(L_i, NL_i, NS_i), i = 1, \dots, 5$, STPR, and TNR in Example 3

	$(n, p) = (500, 400)$					STPR	TNR
	(L_1, NL_1, NS_1)	(L_2, NL_2, NS_2)	(L_3, NL_3, NS_3)	(L_4, NL_4, NS_4)	(L_5, NL_5, NS_5)		
AWG-Lasso+HDIC	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(1.0, 0.0, 0.0)	(1.0, 0.0, 0.0)	(1.0, 0.0, 0.0)	1.0	0.997
AWG-Lasso+HDIC $_{\Pi}$	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(1.0, 0.0, 0.0)	(1.0, 0.0, 0.0)	(1.0, 0.0, 0.0)	1.0	0.997
Rqpen	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.48, 0.52, 0.0)	(0.40, 0.60, 0.0)	(0.42, 0.58, 0.0)	0.66	0.406
K-method	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	0.6	1.0

5 Proofs of the main theorems

First we introduce notation and assumptions. Then we prove Theorems 1 and 2. All the technical proofs are given in the supplement. We denote the conditional probability and expectation on $\{(\mathbf{X}_i, Z_i)\}_{i=1}^n$ by $P_\epsilon(\cdot)$ and $E_\epsilon(\cdot)$, respectively.

Assumption A1 is about $|\mathcal{S}_c^0|$ and $|\mathcal{S}_v^0|$.

Assumption A1: There are bounded constants C_c, C_v, M_c , and M_v such that $|\mathcal{S}_c^0| \leq C_c < M_c$ and $|\mathcal{S}_v^0| \leq C_v < M_v$. Besides, we know M_c and M_v in advance.

This assumption looks restrictive and we may be able to relax this assumption slightly. However, there are still many assumptions and parameters and we decided not to introduce more complications to relax Assumption A1. Note that we can easily relax the conditions on C_c only for Theorem 1 if $\sum_{j \in \mathcal{S}_c^0} w_{1j}^2 = O_p(1)$.

Assumptions A2 and A2' are about the relevant non-zero coefficients and coefficient functions. We need to assume that they are large enough to be detected for our consistency results. Recall that L is the dimension of the spline basis and referred to in Assumption A3 and that q_n appeared in (7).

Assumption A2: We have in probability

$$\frac{\min_{j \in \mathcal{S}_c^0} |g_{cj}| \wedge \min_{j \in \mathcal{S}_v^0} \|g_{vj}\|}{L^{1/2} \{(n^{-1} \log p_n)^{1/2} + \lambda |w_{\mathcal{S}_c^0}|\}} \rightarrow \infty.$$

Assumption A2': We have

$$\frac{\min_{j \in \mathcal{S}_c^0} |g_{cj}| \wedge \min_{j \in \mathcal{S}_v^0} \|g_{vj}\|}{q_n^{1/2} (n^{-1} L \log p_n)^{1/2}} \rightarrow \infty.$$

Next we consider the smoothness of relevant non-zero coefficient functions and spline approximation.

Assumption A3: We take $L = c_L n^{1/5}$ and use linear or smoother splines. Besides, we have for some positive C_g ,

$$\sum_{j \in \mathcal{S}_c^0 \cup \mathcal{S}_v^0} (\|g_j\|_\infty + \|g'_j\|_\infty + \|g''_j\|_\infty) \leq C_g.$$

When Assumption A3 holds, there exists $\boldsymbol{\gamma}_j^* = (\gamma_{1j}^*, \boldsymbol{\gamma}_{-1j}^{*T})^T \in R^L$ for every $j \in \mathcal{S}_c^0 \cup \mathcal{S}_v^0$ such that

$$\sum_{j \in \mathcal{S}_c^0 \cup \mathcal{S}_v^0} \|g_j - \boldsymbol{\gamma}_j^{*T} \mathbf{B}\|_\infty \leq C_1 L^{-2}, \quad \gamma_{1j}^* = L^{1/2} g_{cj}, \quad \text{and} \quad \sum_{j \in \mathcal{S}_v^0} \|g_{vj} - \boldsymbol{\gamma}_{-1j}^{*T} \mathbf{B}_{-1}\|_\infty \leq C_2 L^{-2},$$

where C_1 and C_2 depend only on C_g and the order of the spline basis. Let $\boldsymbol{\gamma}_{\mathcal{S}^0}^*$ consist of γ_{1j}^* , $j \in \mathcal{S}_c^0$, and $\boldsymbol{\gamma}_{-1j}^*$, $j \in \mathcal{S}_v^0$. For \mathcal{S} including the true \mathcal{S}^0 , $\boldsymbol{\gamma}_{\mathcal{S}}^*$ means a vector of coefficients for our spline basis to approximate g_j up to the order of L^{-2} . When $j \in \mathcal{S}_c \cap \overline{\mathcal{S}_c^0}$ or $j \in \mathcal{S}_v \cap \overline{\mathcal{S}_v^0}$, the corresponding elements are put to 0. The other elements are γ_{1j}^* , $j \in \mathcal{S}_c^0$, and $\boldsymbol{\gamma}_{-1j}^*$, $j \in \mathcal{S}_v^0$. See Section S.3 in the supplement for more details on the above approximations.

We define some notation related to spline approximation, δ_i , δ_{ij} , ϵ'_i , and τ_i , by $\delta_{ij} = g_j(Z_j) - \boldsymbol{\gamma}_j^{*T} \mathbf{B}(Z_i)$,

$$\begin{aligned} \delta_i &= \sum_{j \in \mathcal{S}_c^0 \cup \mathcal{S}_v^0} X_{ij} (g_j(Z_i) - \boldsymbol{\gamma}_j^{*T} \mathbf{B}(Z_i)) = \sum_{j \in \mathcal{S}_c^0 \cup \mathcal{S}_v^0} X_{ij} \delta_{ij}, \\ \epsilon'_i &= \epsilon_i + \delta_i, \quad \text{and} \quad \tau_i = \mathbb{P}_\epsilon(\epsilon'_i \leq 0). \end{aligned} \quad (22)$$

Under Assumptions A3 and A4 below, we have uniformly in i and j ,

$$|\delta_{ij}| = O(L^{-2}) \quad \text{and} \quad |\delta_i| \leq C_1 X_M L^{-2} \rightarrow 0$$

for some positive C_1 , where let X_M be a constant satisfying $\max_{i,j} |X_{ij}| \leq X_M$. We allow X_M to diverge as in Assumptions A4 and A4'. Note that

$$\frac{1}{n} \sum_{i=1}^n \delta_i^2 \leq \left\{ n^{-1} \sum_{i=1}^n \left(\sum_{j \in \mathcal{S}_c^0 \cup \mathcal{S}_v^0} X_{ij}^2 \right)^2 \right\}^{1/2} \left\{ n^{-1} \sum_{i=1}^n \left(\sum_{j \in \mathcal{S}_c^0 \cup \mathcal{S}_v^0} \delta_{ij}^2 \right)^2 \right\}^{1/2}. \quad (23)$$

When we examine the properties of our BIC type criteria, we need more smoothness of the coefficient functions to evaluate the approximation bias. We replace Assumption A3 with Assumption A3' for simplicity of presentation. In fact, the Hölder continuity of g_j'' with exponent $\alpha \geq 1/2$ is sufficient if $X_M^4 L^{-2\alpha} = O(L^{-1})$. If X_M is bounded, the proof of Theorem 2 will work for $\alpha = 1/2$. See Lemma 4 in Subsection S.2.2 of the supplement. When we assume Assumption A3', we can replace L^{-2} with L^{-3} in the above approximations.

Assumption A3': We take $L = c_L n^{1/5}$ and use quadratic or smoother splines. Besides, we have for some positive C_g ,

$$\sum_{j \in \mathcal{S}_c^0 \cup \mathcal{S}_v^0} (\|g_j\|_\infty + \|g_j'\|_\infty + \|g_j''\|_\infty + \|g_j^{(3)}\|_\infty) \leq C_g.$$

Next we state assumptions on X_M , p , and q_n . When we consider additive models, we can take $X_M = 1$. Assumptions A4 and A4' imply that ι in $p = O(\exp(n^\iota))$ is less than $1/5$.

Assumption A4: For any positive k ,

$$X_M(\log p_n)^{1/2}n^{-1/10}(\log n)^k \rightarrow 0. \quad (24)$$

Besides, $E\{B_{0l}^2(Z_1)X_{1j}^2\} = O(L^{-1})$ and $E\{B_{0l}(Z_1)|X_{1j}\} = O(L^{-1})$ uniformly in l and j . Recall that $B_{0l}(z)$ is the l -th element of the B-spline basis.

Assumption A4': In Assumption A4, (24) is replaced with

$$X_M(\log p_n)^{1/2}q_n^{3/2}n^{-1/10}(\log n)^k \rightarrow 0.$$

Next we state assumptions on the conditional distribution of ϵ_i on (\mathbf{X}_i, Z_i) . We denote the conditional distribution function by $F_i(\epsilon)$ and the conditional density function by $f_i(\epsilon)$.

Assumption A5: There exist positive C_{f1} , C_{f2} , and C_{f3} such that uniformly in i ,

$$|F_i(u + \delta) - F_i(\delta) - uf_i(\delta)| \leq C_{f1}u^2 \text{ and } f_i(\delta) \leq C_{f2} \text{ when } |\delta| + |u| \leq C_{f3}.$$

Assumption A5': In addition to Assumption A5, $E\{|\epsilon_i|\} < \infty$ and when $|a| \rightarrow 0$, we have uniformly in i ,

$$E_\epsilon[(a - \epsilon_i - \delta_i)I\{0 < \epsilon_i + \delta_i \leq a\}] = \frac{a^2}{2}f_i(-\delta_i) + O(|a|^3) \text{ for } a > 0,$$

and

$$E_\epsilon[(\epsilon_i + \delta_i - a)I\{a < \epsilon_i + \delta_i \leq 0\}] = \frac{a^2}{2}f_i(-\delta_i) + O(|a|^3) \text{ for } a < 0.$$

Actually, when $a > 0$ and $a \rightarrow 0$, we have under some regularity conditions that

$$\int_{-\delta_i}^{a-\delta_i} (a - \epsilon_i - \delta_i)f_i(\epsilon)d\epsilon = \frac{a^2}{2}f_i(-\delta_i) + O(a^3).$$

We introduce some more notation and another kind of assumptions to describe properties of the adaptively weighted Lasso estimators.

We define two index sets \mathbf{S}_M and \mathbf{S}_{C+M} . These index sets are defined for Theorem 2 and they are related to Assumption A1.

$$\mathbf{S}_M = \{\mathcal{S} \mid \mathcal{S}^0 \subset \mathcal{S}, |\mathcal{S}_c| \leq M_c, \text{ and } |\mathcal{S}_v| \leq M_v\} \quad \text{and} \quad (25)$$

$$\mathbf{S}_{C+M} = \{\mathcal{S} \mid \mathcal{S}^0 \subset \mathcal{S}, |\mathcal{S}_c| \leq C_c + M_c, \text{ and } |\mathcal{S}_v| \leq C_v + M_v\} \quad (26)$$

We define some random variables related to $\mathbf{W}_{i\mathcal{S}}$ and describe assumptions on those random variables. The assumptions on those random variables follow from similar assumptions on their population versions and standard technical arguments. We omit the assumptions on the population versions and standard technical arguments here since they are just standard ones in the literature.

We define $\Theta_1(\mathcal{S})$ by

$$\Theta_1(\mathcal{S}) = \frac{1}{n} \sum_{i=1}^n |\mathbf{W}_{i\mathcal{S}}|^2 = \frac{1}{n} \sum_{i=1}^n L^{-1} \sum_{j \in \mathcal{S}_c} |X_{ij}|^2 + \frac{1}{n} \sum_{i=1}^n |\mathbf{B}_{-1}(Z_i)|^2 \sum_{j \in \mathcal{S}_v} |X_{ij}|^2.$$

For technical and notational convenience, we redefine $\Theta_1(\mathcal{S})$ by $\Theta_1(\mathcal{S}) \vee 1$.

Assumption B1: For some positive C_{B1} , we have $\Theta_1(\mathcal{S}^0) \leq C_{B1}$ with probability tending to 1,

Assumption B1 follows from some mild moment conditions under Assumption A1.

We define $\Theta_2(\mathcal{S})$ and $\Theta_3(\mathcal{S})$ by

$$\Theta_2(\mathcal{S}) = L\lambda_{\min}(\widehat{\Sigma}_{\mathcal{S}}) \quad \text{and} \quad \Theta_3(\mathcal{S}) = L\lambda_{\max}(\widehat{\Sigma}_{\mathcal{S}}),$$

where $\widehat{\Sigma}_{\mathcal{S}} = n^{-1} \sum_{i=1}^n f_i(-\delta_i) \mathbf{W}_{i\mathcal{S}} \mathbf{W}_{i\mathcal{S}}^T$. The following assumptions are about their eigenvalues. Recall that our normalization factor of the basis is L^{-1} .

Assumption B2: For some positive C_{B2} , we have $\Theta_2(\mathcal{S}^0) \geq C_{B2}$ with probability tending to 1.

Assumption B2': For some positive C'_{B2} , we have $\Theta_2(\mathcal{S}) \geq C'_{B2}$ uniformly in $\mathcal{S} \in \mathbf{S}_{C+M}$ with probability tending to 1.

Assumption B3: For some positive C_{B3} , we have with probability tending to 1

$$\begin{aligned} \Theta_3(\mathcal{S}^0 \cup (\{j\}, \phi)) &\leq C_{B3} \quad \text{uniformly in } j \in \overline{\mathcal{S}_c^0} \quad \text{and} \\ \Theta_3(\mathcal{S}^0 \cup (\phi, \{j\})) &\leq C_{B3} \quad \text{uniformly in } j \in \overline{\mathcal{S}_v^0}. \end{aligned}$$

Assumption B3': For some positive C'_{B3} , we have with probability tending to 1

$$\Theta_3(\mathcal{S}) \leq C'_{B3} \quad \text{uniformly in } \mathcal{S} \in \mathbf{S}_{C+M}.$$

We define Θ_4 by $\Theta_4 = n^{-1} \sum_{i=1}^n \sum_{j \in \mathcal{S}_v^0} X_{ij}^2$.

Assumption B4: For some positive C_{B4} , we have $\Theta_4 \leq C_{B4}$ with probability tending to 1.

Assumption B4': In addition to Assumption B4, we have for some positive C'_{B4} ,

$$n^{-1} \sum_{i=1}^n \left(\sum_{j \in \mathcal{S}_c^0 \cup \mathcal{S}_v^0} X_{ij}^2 \right)^2 \leq C'_{B4} \quad \text{with probability tending to 1.}$$

Assumption B4' is used to control (23). Assumptions B4 and B4' follow from mild moment conditions under Assumption A1.

We define $\Theta_5(\mathcal{S})$ by $\Theta_5(\mathcal{S}) = \max_{1 \leq i \leq n} |\mathbf{W}_{i\mathcal{S}}|^2$. Notice that there are positive constants C_1 and C_2 such that

$$|\mathbf{W}_{i\mathcal{S}}|^2 = L^{-1} \sum_{j \in \mathcal{S}_c} X_{ij}^2 + |\mathbf{B}_{-1}(Z_i)|^2 \sum_{j \in \mathcal{S}_v} X_{ij}^2 \leq C_1 X_M^2 (L^{-1} |\mathcal{S}_c| + |\mathcal{S}_v|) \leq C_2 X_M^2 \quad (27)$$

for any $\mathcal{S} \in \mathcal{S}_{C+M}$ under Assumption A1.

We define $\widehat{\Omega}_{\mathcal{S}}$ by $\widehat{\Omega}_{\mathcal{S}} = n^{-1} \sum_{i=1}^n \tau_i (1 - \tau_i) \mathbf{W}_{i\mathcal{S}} \mathbf{W}_{i\mathcal{S}}^T$. The last assumption is about its eigenvalues. Recall that τ_i is defined in (22).

Assumption B5: There is a positive constant C_{B5} such that uniformly in $\mathcal{S} \in \mathcal{S}_{C+M}$,

$$\frac{1}{C_{B5}} \leq L \lambda_{\min}(\widehat{\Omega}_{\mathcal{S}}) \leq L \lambda_{\max}(\widehat{\Omega}_{\mathcal{S}}) \leq C_{B5} \quad \text{with probability tending to 1.}$$

We state Proposition 1 before we prove Theorem 1. The proposition gives the convergence rate of the AWG-Lasso estimator. We prove this proposition by following that of Theorem 1 in [7] in the supplement.

We use the proposition with $\mathcal{S} = \mathcal{S}^0$ or with $\mathcal{S} \in \mathcal{S}_{C+M}$ and $\lambda = 0$. Let $w_{\mathcal{S}}$ be a vector consisting of $\{w_{1j} \mid j \in \mathcal{S}_c\}$ and $\{w_{-1j} \mid j \in \mathcal{S}_v\}$. Then we define $|w_{\mathcal{S}}|$ and K_n by

$$|w_{\mathcal{S}}|^2 = \sum_{j \in \mathcal{S}_c} w_{1j}^2 + \sum_{j \in \mathcal{S}_v} w_{-1j}^2 \quad \text{and} \quad K_n(\mathcal{S}) = \sqrt{n^{-1} \Theta_1(\mathcal{S}) \log p_n} + \lambda |w_{\mathcal{S}}|.$$

Tentatively we assume the weights are constants, not random variables.

Proposition 1 *Suppose that $\mathcal{S}^0 \subset \mathcal{S}$ and Assumptions A1 and A3-5 hold. Besides we assume*

$$\left(\frac{\Theta_5(\mathcal{S})}{\Theta_2(\mathcal{S})} \right)^{1/2} (\Theta_2^{-1/2}(\mathcal{S}) \vee \Theta_4^{1/2}) K_n(\mathcal{S}) L \rightarrow 0 \quad (28)$$

and we define η_n by $\eta_n = C_M L K_n(\mathcal{S})$, where C_M satisfies

$$C_M \geq b_1 \left\{ \frac{1}{\Theta_2(\mathcal{S})} \vee \left(\frac{\Theta_4}{\Theta_2(\mathcal{S})} \right)^{1/2} \right\} \quad (29)$$

for sufficiently large b_1 depending on b_2 in (30). Then we have for any fixed positive b_2 that

$$P_\epsilon(|\widehat{\gamma}_S^\lambda - \gamma_S^*| \geq \eta_n) \leq \exp(-b_2 \log p_n). \quad (30)$$

Later we use Assumptions B1-4 to control random variables in (28) and (29) in Proposition 1. Here some remarks on Proposition 1 are in order.

Remark 1 When w_S is a random vector and $\lambda > 0$, “ $\rightarrow 0$ ” in (28) should be replaced with “ $\xrightarrow{P} 0$.” Besides, when for some positive C_1 , C_2 , and C_3 ,

$$P(C_1 \leq \Theta_2(\mathcal{S}), \Theta_1(\mathcal{S}) \leq C_2, \Theta_4 \leq C_3) \rightarrow 1,$$

the RHS of (29) is bounded from above in probability and $\Theta_1(\mathcal{S})$ in $K_n(\mathcal{S})$ can be replaced with a constant. Thus we have $P(|\widehat{\gamma}_S^\lambda - \gamma_S^*| \geq \eta_n) \rightarrow 0$ under (28) in probability with a fixed C_M . Especially when $\mathcal{S} = \mathcal{S}^0$,

$$\eta_n \sim L\{(n^{-1} \log p_n)^{1/2} + \lambda|w_{\mathcal{S}^0}|\}.$$

Remark 2 Since $\Theta_5(\mathcal{S}^0) \leq C_4 X_M^2$ for some positive C_4 under Assumption A1, (28) reduces to $X_M L\{(n^{-1} \log p_n)^{1/2} + \lambda|w_{\mathcal{S}^0}|\} \xrightarrow{P} 0$ in the setup of Remark 1 with $\mathcal{S} = \mathcal{S}^0$ and this is not a restrictive condition.

Remark 3 When $\lambda = 0$ and the assumptions in Theorem 2 hold, we have for $\widehat{\gamma}_S^\lambda = \widetilde{\gamma}_S$ that

$$|\widehat{\gamma}_S^\lambda - \gamma_S^*| = |\widetilde{\gamma}_S - \gamma_S^*| \leq C_5 L(n^{-1} \log p_n)^{1/2}$$

uniformly in $\mathcal{S} \in \mathcal{S}_{C+M}$ with probability tending to 1 for some positive C_5 . We use this result in the proof of Theorem 2.

We provide the proof of Theorem 1. We define $\Gamma_S(M)$ by

$$\Gamma_S(M) = \{\gamma_S \in R^{d_V(\mathcal{S})} \mid |\gamma_S - \gamma_S^*| \leq M\} \quad (31)$$

Proof of Theorem 1 First we prove $(\widehat{\gamma}_{\mathcal{S}^0}^\lambda, \mathbf{0}^T)^T \in R^{pL}$ is a global minimizer of (5) by checking the following conditions (32) and (33). These conditions follow from the standard optimization theory as in [38] and [28]. In addition to (32) as in [38] and [28], we should deal with (33) since we are employing group penalties. Hereafter in this proof, we omit the superscript λ and write $\widehat{\gamma}_{\mathcal{S}^0}$ for $\widehat{\gamma}_{\mathcal{S}^0}^\lambda$

With probability tending to 1, we have

$$\left| \frac{1}{n} \sum_{i=1}^n L^{-1/2} X_{ij} \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \widehat{\gamma}_{S^0}) \right| \leq \lambda w_{1j} \text{ for any } j \in \overline{\mathcal{S}}_c^0 \quad \text{and} \quad (32)$$

$$\left| \frac{1}{n} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \widehat{\gamma}_{S^0}) \right| \leq \lambda w_{-1j} \text{ for any } j \in \overline{\mathcal{S}}_v^0. \quad (33)$$

We verify only (33) since (32) is easier.

Proposition 1, Remark 1, and the conditions of this theorem imply that

$$|\widehat{\gamma}_{S^0} - \gamma_{S^0}^*| \leq C_1 L \{(n^{-1} \log p_n)^{1/2} + \lambda |w_{S^0}|\} \leq C_2 L (n^{-1} \log p_n)^{1/2} (\log n)^{k_\lambda} \quad (34)$$

with probability tending to 1 for some positive C_1 and C_2 . We define $V_j(\gamma_{S^0})$ by

$$\begin{aligned} V_j(\gamma_{S^0}) &= n^{-1} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} \left\{ \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}) - \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}^*) \right\} \\ &\quad - \mathbb{E}_\epsilon \left[n^{-1} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} \left\{ \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}) - \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}^*) \right\} \right] \end{aligned}$$

By considering the upper bounds given in (34), we can take a positive constant C_ξ for any small positive ξ such that with probability larger than $1 - \xi$,

$$\begin{aligned} &\left| \frac{1}{n} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \widehat{\gamma}_{S^0}) \right| \quad (35) \\ &\leq \left| \mathbb{E}_\epsilon \left[\frac{1}{n} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} \left\{ \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}) - \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}^*) \right\} \right]_{\gamma_{S^0} = \widehat{\gamma}_{S^0}} \right| \\ &+ \left| \frac{1}{n} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}^*) \right| + \max_{\gamma_{S^0} \in \Gamma_{S^0}(C_\xi L (n^{-1} \log p_n)^{1/2} (\log n)^{k_\lambda})} |V_j(\gamma_{S^0})|. \end{aligned}$$

We use the following two lemmas to evaluate (35). These lemmas are to be proved in the supplement.

Lemma 1 *For some positive C_1 , we have*

$$\left| \frac{1}{n} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}^*) \right| \leq C_1 (n^{-1} \log p_n)^{1/2}$$

uniformly in $j \in \overline{\mathcal{S}}_v^0$ with probability tending to 1

Lemma 2 *Take any fixed positive C and k and fix them. Then we have*

$$\max_{\gamma_{S^0} \in \Gamma_{S^0}(CL(n^{-1} \log p_n)^{1/2} (\log n)^k)} |V_j(\gamma_{S^0})| = o_p(\lambda)$$

uniformly in $j \in \overline{\mathcal{S}}_v^0$.

Finally we evaluate

$$\begin{aligned} & \mathbb{E}_\epsilon \left[\frac{1}{n} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} \{ \rho'_\tau(Y_i - \mathbf{W}_{i\mathcal{S}^0}^T \gamma_{\mathcal{S}^0}) - \rho'_\tau(Y_i - \mathbf{W}_{i\mathcal{S}^0}^T \gamma_{\mathcal{S}^0}^*) \} \right]_{\gamma_{\mathcal{S}^0} = \hat{\gamma}_{\mathcal{S}^0}} \quad (36) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} \{ F_i(-\delta_i) - F_i(-\delta_i + \mathbf{W}_{i\mathcal{S}^0}^T (\hat{\gamma}_{\mathcal{S}^0} - \gamma_{\mathcal{S}^0}^*)) \}. \end{aligned}$$

Setting $\hat{\Delta}^0 = \hat{\gamma}_{\mathcal{S}^0} - \gamma_{\mathcal{S}^0}^*$ and recalling Assumption A5, we find that (36) is rewritten as

$$-\frac{1}{n} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} f_i(-\delta_i) \mathbf{W}_{i\mathcal{S}^0}^T \hat{\Delta}^0 + o_p((n^{-1} \log p_n)^{1/2}) = -D_j \hat{\Delta}^0 + o_p((n^{-1} \log p_n)^{1/2}) \quad (37)$$

uniformly in $j \in \overline{\mathcal{S}_v^0}$, where D_j is clearly defined in the above equation.

Assumption B3 implies that for some positive C_1 ,

$$\lambda_{\max}(D_j^T D_j) \leq C_1 L^{-2} \quad (38)$$

uniformly in $j \in \overline{\mathcal{S}_v^0}$ with probability tending to 1. This is because D_j is part of $\widehat{\Sigma}_{\mathcal{S}^0 \cup \{\phi, \{j\}\}}$. Thus (34) and (38) yield that for some positive C_2 ,

$$|D_j \hat{\Delta}^0| \leq C_2 \{ (n^{-1} \log p_n)^{1/2} + \lambda |w_{\mathcal{S}^0}| \} \quad (39)$$

uniformly in $j \in \overline{\mathcal{S}_v^0}$ with probability tending to 1.

By combining (35), Lemmas 1 and 2, (37), and (39), we obtain

$$\left| \frac{1}{n} \sum_{i=1}^n \mathbf{B}_{-1}(Z_i) X_{ij} \rho'_\tau(Y_i - \mathbf{W}_{i\mathcal{S}^0}^T \hat{\gamma}_{\mathcal{S}^0}) \right| \leq \lambda w_{-1j}$$

uniformly in $j \in \overline{\mathcal{S}_v^0}$ with probability tending to 1. Hence (33) is established.

As for the latter part of the theorem, Assumption A2 implies that γ_{1j}^* , $j \in \mathcal{S}_c^0$, and γ_{-1j}^* , $j \in \mathcal{S}_v^0$, are large enough to be detected due to Proposition 1 with $\mathcal{S} = \mathcal{S}^0$.

Hence the proof of the theorem is complete.

Now we state the proof of Theorem 2

Proof of Theorem 2) We give the details of the overfitting case here. We can deal with the underfitting case by following the standard arguments and we give the proof of the underfitting case in the supplement.

Let \mathcal{S} satisfy $\mathcal{S} \in \mathbf{S}_M$ and $\mathcal{S} \neq \mathcal{S}^0$. See (25) for the definition of \mathbf{S}_M . “Uniformly in \mathcal{S} ” means “uniformly in \mathcal{S} satisfying $\mathcal{S} \in \mathbf{S}_M$ and $\mathcal{S} \neq \mathcal{S}^0$ ”. We have replaced Assumption A3 with Assumption A3’. We use Assumption A3’ only once in the proof

(Lemma 4) and we use Assumption A3 in the other part. Assumption A3' can be relaxed in some cases. See Lemma 4 in Subsection S.2.2 of the supplement for more details.

If we have established

$$R_V(\boldsymbol{\gamma}_{S^0}^*) = \frac{1}{n} \sum_{i=1}^n \rho_\tau(\epsilon_i) + O(X_M L^{-2}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}\{\rho_\tau(\epsilon_i)\} + o_p(1), \quad (40)$$

$$R_V(\tilde{\boldsymbol{\gamma}}_{S^0}) = R_V(\boldsymbol{\gamma}_{S^0}^*) + o_p(1), \quad \text{and uniformly in } \mathcal{S}, \quad (41)$$

$$R_V(\tilde{\boldsymbol{\gamma}}_{S^0}) - R_V(\tilde{\boldsymbol{\gamma}}_S) = (d_V(\mathcal{S}) - d_V(\mathcal{S}^0))O_p(n^{-1}\{(\log p_n) \vee (q_n \log p_n)^{1/2}\}), \quad (42)$$

then we have for some positive C_1 ,

$$\begin{aligned} 0 &\leq \log R_V(\tilde{\boldsymbol{\gamma}}_{S^0}) - \log R_V(\tilde{\boldsymbol{\gamma}}_S) = -\log \left\{ 1 + \frac{R_V(\tilde{\boldsymbol{\gamma}}_S) - R_V(\tilde{\boldsymbol{\gamma}}_{S^0})}{R_V(\tilde{\boldsymbol{\gamma}}_{S^0})} \right\} \\ &\leq \frac{1}{C_1} \{R_V(\tilde{\boldsymbol{\gamma}}_{S^0}) - R_V(\tilde{\boldsymbol{\gamma}}_S)\} \end{aligned} \quad (43)$$

uniformly in \mathcal{S} with probability tending to 1. By (42) and (43), we obtain

$$\begin{aligned} \log R_V(\tilde{\boldsymbol{\gamma}}_{S^0}) - \log R_V(\tilde{\boldsymbol{\gamma}}_S) &= (d_V(\mathcal{S}) - d_V(\mathcal{S}^0))O_p(n^{-1}\{\log p_n \vee (q_n \log p_n)^{1/2}\}) \\ &< (d_V(\mathcal{S}) - d_V(\mathcal{S}^0))\frac{\log p_n}{2n}q_n \end{aligned}$$

uniformly in \mathcal{S} with probability tending to 1. Hence the proof for the overfitting case is complete.

Thus we have only to prove (40)-(42). We prove only (42) since (40) and (41) are easy to deal with.

(49), (50), and (53), which will be defined later, are important when we prove (42). To verify (49), first we will prove in the supplement that

$$\begin{aligned} R_V(\boldsymbol{\gamma}_S) - R_V(\boldsymbol{\gamma}_S^*) &= -(\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*)^T \frac{1}{n} \sum_{i=1}^n \mathbf{W}_{iS}(\tau_i - I\{\epsilon'_i \leq 0\}) + \frac{1}{2}(\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*)^T \widehat{\Sigma}_S(\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*) \\ &\quad + (\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*)^T \frac{1}{n} \sum_{i=1}^n \mathbf{W}_{iS}(\tau_i - \tau) + O_p\left(\frac{\log p_n}{n(\log n)^2}\right) \end{aligned} \quad (44)$$

uniformly in $\boldsymbol{\gamma}_S \in \Gamma_S(M_1 L(q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} for any fixed M_1 .

We use (44) to derive a useful expression of $R_V(\tilde{\boldsymbol{\gamma}}_S)$. Put

$$\mathbf{a}_S = \frac{1}{n} \sum_{i=1}^n \mathbf{W}_{iS}(\tau_i - I\{\epsilon'_i \leq 0\}), \quad \mathbf{b}_S = \frac{1}{n} \sum_{i=1}^n \mathbf{W}_{iS}(\tau_i - \tau), \quad \text{and } \bar{\boldsymbol{\gamma}}_S - \boldsymbol{\gamma}_S^* = \widehat{\Sigma}_S^{-1} \mathbf{a}_S. \quad (45)$$

According to (S.19) in Lemma 4 in Subsection S.2.2 of the supplement,

$$(\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*)^T \frac{1}{n} \sum_{i=1}^n \mathbf{W}_{iS}(\tau_i - \tau) = (\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*)^T \mathbf{b}_S = O_p\left(\frac{(q_n \log p_n)^{1/2}}{n}\right) \quad (46)$$

and this term in (44) is negligible uniformly in $\gamma_S \in \Gamma_S(M_1 L(q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} for any fixed M_1 .

By applying Bernstein's inequality conditionally on $\{(\mathbf{X}_i, Z_i)\}_{i=1}^n$ first and using Assumption B5, we have

$$|\mathbf{a}_S|^2 = O_p\left(\frac{\log p_n}{n}\right) \quad (47)$$

uniformly in \mathcal{S} . Thus we have from Assumption B2' that uniformly in \mathcal{S} ,

$$\bar{\gamma}_S - \gamma_S^* = O_p(L(n^{-1} \log p_n)^{1/2}). \quad (48)$$

We take some $\boldsymbol{\delta}_S \in R^{d_V(\mathcal{S})}$. If $\bar{\gamma}_S + \boldsymbol{\delta}_S \in \Gamma_S(M_1 L(q_n n^{-1} \log p_n)^{1/2})$, we have from (44) and (46) that uniformly in $\boldsymbol{\delta}_S$ and \mathcal{S} ,

$$R_V(\bar{\gamma}_S + \boldsymbol{\delta}_S) - R_V(\gamma_S^*) = -\frac{1}{2} \mathbf{a}_S^T \widehat{\Sigma}_S^{-1} \mathbf{a}_S + \frac{1}{2} \boldsymbol{\delta}_S^T \widehat{\Sigma}_S \boldsymbol{\delta}_S + O_p\left(\frac{(q_n \log p_n)^{1/2}}{n}\right) + O_p\left(\frac{\log p_n}{n(\log n)^2}\right). \quad (49)$$

Because of the optimality of $R_V(\tilde{\gamma}_S)$ and (49), we should have

$$R_V(\tilde{\gamma}_S) - R_V(\gamma_S^*) = -\frac{1}{2} \mathbf{a}_S^T \widehat{\Sigma}_S^{-1} \mathbf{a}_S + O_p\left(\frac{(q_n \log p_n)^{1/2}}{n}\right) + O_p\left(\frac{\log p_n}{n(\log n)^2}\right) \quad (50)$$

uniformly in \mathcal{S} . The above arguments show that this expression also holds for \mathcal{S}^0 . By combining (49) and (50) and setting $\boldsymbol{\delta}_S = \tilde{\gamma}_S - \bar{\gamma}_S$, we also obtain

$$|\tilde{\gamma}_S - \bar{\gamma}_S|^2 = O_p\left(\frac{L(q_n \log p_n)^{1/2}}{n}\right) + O_p\left(\frac{L \log p_n}{n(\log n)^2}\right) \quad (51)$$

uniformly in \mathcal{S} . Note again that these expressions also hold for \mathcal{S}^0 . This equation is used later in the underfitting case.

We evaluate the difference between $R_V(\tilde{\gamma}_S)$ and $R_V(\tilde{\gamma}_{S^0})$. Now write

$$\widehat{\Sigma}_S = \begin{pmatrix} \widehat{\Sigma}_{S^0} & \widehat{\Sigma}_{S12} \\ \widehat{\Sigma}_{S21} & \widehat{\Sigma}_{S22} \end{pmatrix} \quad \text{and} \quad \mathbf{a}_S = \begin{pmatrix} \mathbf{a}_{S^0} \\ \mathbf{a}_{S2} \end{pmatrix} \quad (52)$$

and notice that $R_V(\gamma_S^*) = R_V(\gamma_{S^0}^*)$. Thus due to (50), we have only to consider the difference

$$\begin{aligned} \mathbf{a}_S^T \widehat{\Sigma}_S^{-1} \mathbf{a}_S - \mathbf{a}_{S^0}^T \widehat{\Sigma}_{S^0}^{-1} \mathbf{a}_{S^0} &= \mathbf{a}_{S^0}^T \widehat{\Sigma}_{S^0}^{-1} \widehat{\Sigma}_{S12} \widehat{F}_{S2} \widehat{\Sigma}_{S21} \widehat{\Sigma}_{S^0}^{-1} \mathbf{a}_{S^0} \\ &\quad - 2 \mathbf{a}_{S^0}^T \widehat{\Sigma}_{S^0}^{-1} \widehat{\Sigma}_{S12} \widehat{F}_{S2} \mathbf{a}_{S2} + \mathbf{a}_{S2}^T \widehat{F}_{S2} \mathbf{a}_{S2}, \end{aligned} \quad (53)$$

where $\widehat{F}_{S2} = (\widehat{\Sigma}_{S22} - \widehat{\Sigma}_{S21} \widehat{\Sigma}_{S^0}^{-1} \widehat{\Sigma}_{S12})^{-1}$, when we evaluate $R_V(\tilde{\gamma}_S) - R_V(\tilde{\gamma}_{S^0})$.

We will demonstrate that the RHS of (53) has the stochastic order of $(d_V(\mathcal{S}) - d_V(\mathcal{S}^0))O_p(n^{-1} \log p_n)$ uniformly in \mathcal{S} .

From Assumptions B2' and B3', we have for some positive C_1 , C_2 , and C_3 ,

$$C_1 L \leq \lambda_{\min}(\widehat{F}_{S2}) \leq \lambda_{\max}(\widehat{F}_{S2}) \leq C_2 L \text{ and } \lambda_{\max}(\widehat{\Sigma}_{S21} \widehat{\Sigma}_{S12}) \leq C_3 L^{-2} \quad (54)$$

uniformly in \mathcal{S} with probability tending to 1.

By applying Bernstein's inequality conditionally on $\{(\mathbf{X}_i, Z_i)\}_{i=1}^n$ first and using Assumption B5, we have that uniformly in \mathcal{S} ,

$$|\mathbf{a}_{S2}|^2 = (d_V(\mathcal{S}) - d_V(\mathcal{S}^0))O_p\left(\frac{\log p_n}{nL}\right). \quad (55)$$

Hence (54) and (55) imply that the third term on the RHS of (53) satisfies

$$\mathbf{a}_{S2}^T \widehat{F}_{S2} \mathbf{a}_{S2} = (d_V(\mathcal{S}) - d_V(\mathcal{S}^0))O_p(n^{-1} \log p_n) \text{ uniformly in } \mathcal{S}. \quad (56)$$

To evaluate the first and second terms on the RHS of (53),

$$(\mathbf{a}_{S0}^T \widehat{\Sigma}_{S0}^{-1} \widehat{\Sigma}_{S12}) \widehat{F}_{S2} (\widehat{\Sigma}_{S21} \widehat{\Sigma}_{S0}^{-1} \mathbf{a}_{S0}) \text{ and } (\mathbf{a}_{S0}^T \widehat{\Sigma}_{S0}^{-1} \widehat{\Sigma}_{S12}) \widehat{F}_{S2} \mathbf{a}_{S2}, \quad (57)$$

we consider

$$\widehat{\Sigma}_{S21} \widehat{\Sigma}_{S0}^{-1} \mathbf{a}_{S0} = \widehat{\Sigma}_{S21} \widehat{\Sigma}_{S0}^{-1} \frac{1}{n} \sum_{i=1}^n \mathbf{W}_{iS0} (\tau_i - I\{\epsilon'_i \leq 0\}) \quad (58)$$

to obtain (62) below. And write

$$\widehat{\Sigma}_{S12} = (\mathbf{s}_1, \dots, \mathbf{s}_{d_V(\mathcal{S}) - d_V(\mathcal{S}^0)})$$

and note that (54) implies

$$\mathbf{s}_j^T \mathbf{s}_j = O_p(L^{-2}) \text{ and } \lambda_{\max}(\widehat{\Sigma}_{S21} \widehat{\Sigma}_{S0}^{-1} \widehat{\Omega}_{S0} \widehat{\Sigma}_{S0}^{-1} \widehat{\Sigma}_{S12}) = O_p(L^{-1}) \quad (59)$$

uniformly in j and \mathcal{S} with probability tending to 1. Besides, we have for some positive C_4 and C_5 ,

$$\max_j |\mathbf{s}_j^T \widehat{\Sigma}_{S0}^{-1} \mathbf{W}_{iS0}| \leq C_4 L |\mathbf{s}_j| |\mathbf{W}_{iS0}| \leq C_5 L |\mathbf{s}_j| X_M = O_p(X_M) \quad (60)$$

uniformly in i and \mathcal{S} with probability tending to 1.

Hence by applying Bernstein's inequality conditionally together with (59) and (60), we obtain

$$\frac{1}{n} \sum_{i=1}^n \mathbf{s}_j^T \widehat{\Sigma}_{S0}^{-1} \mathbf{W}_{iS0} (\tau_i - I\{\epsilon'_i \leq 0\}) = O_p(\{(nL)^{-1} \log p_n\}^{1/2}) \quad (61)$$

uniformly in j and \mathcal{S} . Therefore (61) yields that uniformly in \mathcal{S} ,

$$|\widehat{\Sigma}_{\mathcal{S}21}\widehat{\Sigma}_{\mathcal{S}^0}^{-1}\mathbf{a}_{\mathcal{S}^0}|^2 = (d_V(\mathcal{S}) - d_V(\mathcal{S}^0))O_p((nL)^{-1}\log p_n). \quad (62)$$

Thus (54), (55), (57), and (62) imply that the first and second terms on the RHS of (53) have the stochastic order of $(d_V(\mathcal{S}) - d_V(\mathcal{S}^0))O_p(n^{-1}\log p_n)$ uniformly in \mathcal{S} as in (56). We have demonstrated that the RHS of (53) has the stochastic order of $(d_V(\mathcal{S}) - d_V(\mathcal{S}^0))O_p(n^{-1}\log p_n)$ uniformly in \mathcal{S} .

Hence (42) follows from (50) and this evaluation of (53) and the proof of the overfitting case is complete. The proof of the underfitting case is given in the supplement. Hence the proof is complete.

Acknowledgements

We would like to thank the Editor, the associate editor, and the reviewer for their constructive comments, which greatly improve the presentation of this paper. We also thank Professors Kengo Kato and Yanlin Tang for providing us with their codes to implement K-method and T-method, and Dr. Wen-Ting Wang for his assistance in code development.

References

- [1] A. Belloni and V. Chernozhukov. l_1 -penalized quantile regression in high-dimensional sparse models. *Ann. Statist.*, 39:82–130, 2011.
- [2] P. J. Bickel, y. Ritov, and A. B. Tsybakov. Simultaneous analysis of lasso and dantzig selector. *Ann. Statist.*, 37:1705–1732, 2009.
- [3] P. Bühlmann and S. van de Geer. *Statistics for High-Dimensional Data: Methods Theory and Applications*. Springer, New York, Dordrecht, Heidelberg, London, 2011.
- [4] Z. Cai and Z. Xiao. Semiparametric quantile regression estimation in dynamic models with partially varying coefficients. *J. Econometrics*, 167:413–425, 2012.
- [5] J. Chen and Z. Chen. Extended bayesian information criteria for model selection with large model spaces. *Biometrika*, 95:759–771, 2008.

- [6] M.-Y. Cheng, T. Honda, J. Li, and H. Peng. Nonparametric independence screening and structure identification for ultra-high dimensional longitudinal data. *Ann. Statist.*, 42:1819–1849, 2014.
- [7] J. Fan, Y. Fan, and E. Barut. Adaptive robust variable selection. *Ann. Statist.*, 42:324–351, 2014.
- [8] J. Fan, Y. Feng, and R. Song. Nonparametric independence screening in sparse ultra-high-dimensional additive models. *J. Amer. Statist. Assoc.*, 106:544–557, 2011.
- [9] J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.*, 96:1348 – 1360, 2001.
- [10] J. Fan, Y. Ma, and W. Dai. Nonparametric independence screening in sparse ultra-high dimensional varying coefficient models. *J. Amer. Statist. Assoc.*, 109:1270–1284, 2014.
- [11] J. Fan and R. Song. Sure independence screening in generalized linear models with np-dimensionality. *Ann. Statist.*, 38:3567–3604, 2010.
- [12] J. Fan, L. Xue, and H. Zou. Strong oracle optimality of folded concave penalized estimation. *Ann. Statist.*, 42:819–849, 2014.
- [13] J. Fan, L. Xue, and H. Zou. Multitask quantile regression under the transnormal model. *J. Amer. Statist. Assoc.*, 111:1726–1735, 2016.
- [14] T. Hastie, R. Tibshirani, and M. Wainwright. *Statistical learning with sparsity*. CRC press, Boca Raton, 2015.
- [15] X. He, L. Wang, and H. G. Hong. Quantile-adaptive model-free variable screening for high-dimensional heterogeneous data. *Ann. Statist.*, 41:342–369, 2013.
- [16] T. Honda and R. Yabe. Variable selection and structure identification for varying coefficient cox models. *J. Multivar. Anal.*, 161:103–122, 2017.
- [17] C.-K. Ing and T. L. Lai. A stepwise regression method and consistent model selection for high-dimensional sparse linear models. *Statistica Sinica*, 22:1473–1513, 2011.

- [18] K. Kato. Group lasso for high dimensional sparse quantile regression models. *arXiv preprint arXiv:1103.1458*, 2011.
- [19] Y. Kim and J.-J. Jeon. Consistent model selection criteria for quadratically supported risks. *Ann. Statist.*, 44:2467–2496, 2016.
- [20] E. R. Lee and E. Mammen. Local linear smoothing for sparse high dimensional varying coefficient models. *Electronic Journal of Statistics*, 10:855–894, 2016.
- [21] E. R. Lee, H. Noh, and B. U. Park. Model selection via bayesian information criterion for quantile regression models. *J. Amer. Statist. Assoc.*, 109:216–229, 2014.
- [22] H. Lian. Semiparametric estimation of additive quantile regression models by two-fold penalty. *Journal of Business & Economic Statistics*, 30:337–350, 2012.
- [23] H. Lian, H. Liang, and D. Ruppert. Separation of covariates into nonparametric and parametric parts in high-dimensional partially linear additive models. *Statist. Sinica*, 25:591–607, 2015.
- [24] S. Lv, H. Lin, H. Lian, and J. Huang. Oracle inequalities for sparse additive quantile regression in reproducing kernel hilbert space. *Forthcoming in Ann. Statist.*, 2017.
- [25] S. Ma and X. He. Inference for single-index quantile regression models with profile optimization. *Ann. Statist.*, 44:1234–1268, 2016.
- [26] L. Meier, S. van de Geer, and P. Bühlmann. The group lasso for logistic regression. *J. Royal Statist. Soc. Ser. B*, 70:53–71, 2008.
- [27] L. L. Schumaker. *Spline Functions: Basic Theory 3rd ed.* Cambridge University Press, Cambridge, 2007.
- [28] B. Sherwood and L. Wang. Partially linear additive quantile regression in ultra-high dimension. *Ann. Statist.*, 44:288–317, 2016.
- [29] Y. Tang, X. Song, H. J. Wang, and Z. Zhu. Variable selection in high-dimensional quantile varying coefficient models. *J. Multivar. Anal.*, 122:115–132, 2013.
- [30] R. J. Tibshirani. Regression shrinkage and selection via the lasso. *J. Royal Statist. Soc. Ser. B*, 58:267–288, 1996.

- [31] S. van de Geer. *Empirical Processes in M-estimation*. Cambridge University Press, Cambridge, 2000.
- [32] S. van de Geer. *Estimation and testing under sparsity*. Springer, Switzerland, 2016.
- [33] H. Wang. Forward regression for ultra-high dimensional variable screening. *J. Amer. Statist. Assoc.*, 104:1512–1524, 2009.
- [34] H. J. Wang, Z. Zhu, and J. Zhou. Quantile regression in partially linear varying coefficient models. *Ann. Statist.*, 37:3841–3866, 2009.
- [35] J. Yan and J. Huang. Model selection for cox models with time-varying coefficients. *Biometrics*, 68:419–428, 2012.
- [36] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *J. Royal Statist. Soc. Ser. B*, 68:49–67, 2006.
- [37] H. H. Zhang, G. Cheng, and Y. Liu. Linear or nonlinear? automatic structure discovery for partially linear models. *J. Amer. Statist. Assoc.*, 106:1099–1112, 2011.
- [38] Q. Zheng, L. Peng, and X. He. Globally adaptive quantile regression with ultra-high dimensional data. *Ann. Statist.*, 43:2225–2258, 2015.
- [39] L. Zhu, M. Huang, and R. Li. Semiparametric quantile regression with high-dimensional covariates. *Statist. Sinica*, 22:1379–1401, 2012.
- [40] H. Zou. The adaptive lasso and its oracle properties. *J. Amer. Statist. Assoc.*, 101:1418–1429, 2006.

Supplement to “Adaptively weighted group Lasso for semiparametric quantile regression model”

by Toshio Honda, Ching-Kang Ing, and Wei-Ying Wu

S.1 Additive models

We can deal with additive models in the same way because of the similarity between (4) and (S.2). We describe the specific procedure for additive models in this section. Recall we assume some initial estimates are available here, too.

We have no index variable and assume the additivity and $X_{ij} \in [0, 1]$ for $j = 1, \dots, p$. Hence our model is

$$Y_i = \mu + \sum_{j=1}^p g_j(X_{ij}) + \epsilon_i, \quad (\text{S.1})$$

where $X_{ij} \in [0, 1]$, $\int_0^1 g_j(x)dx = 0$, and $E\{\rho'_\tau(\epsilon_i) | \mathbf{X}_i\} = 0$. To deal with partially linear additive coefficient models, we decompose $g_j(x)$ as $g_j(x) = g_{lj}(x) + g_{aj}(x)$, where $g_{lj}(x) = c_{lj}B_2(x)$ (the j -th linear component) and $g_{aj}(x)$ (the j -th nonlinear component) satisfies

$$\int_0^1 g_{lj}(x)g_{aj}(x)dx = 0.$$

Our regression spline model is given by

$$Y_i = \mu + \mathbf{W}_i^T \boldsymbol{\gamma}_{-1} + \epsilon'_i, \quad (\text{S.2})$$

where $\boldsymbol{\gamma}_{-1} = (\boldsymbol{\gamma}_{-11}^T, \dots, \boldsymbol{\gamma}_{-1p}^T)^T$ and $\mathbf{W}_i = (\mathbf{B}_{-1}^T(X_{i1}), \dots, \mathbf{B}_{-1}^T(X_{ip}))^T$, with $\boldsymbol{\gamma}_{-1j}$ and $\mathbf{B}_{-1}(z)$ defined as in Section 2. Denote the true model by $\mathcal{S}^0 = (\mathcal{S}_l^0, \mathcal{S}_a^0)$, where

$$\mathcal{S}_l^0 = \{j | g_{lj}(x) \not\equiv 0\} \quad \text{and} \quad \mathcal{S}_a^0 = \{j | g_{aj}(x) \not\equiv 0\}.$$

When some j 's satisfy both $j \in \mathcal{S}_l^0$ and $j \notin \mathcal{S}_a^0$ simultaneously, our true model is a partially linear additive model.

We describe the details of our simultaneous variable selection and structure identification procedure for additive models. First express $\boldsymbol{\gamma}_{-1j}$ as $\boldsymbol{\gamma}_{-1j} = (\boldsymbol{\gamma}_{2j}, \boldsymbol{\gamma}_{-2j}^T)^T$, noting that $\boldsymbol{\gamma}_{2j}$ is for $B_2(X_{ij}) = \sqrt{12/L}(X_{ij}-1/2)$ and $\boldsymbol{\gamma}_{-2j}$ is for $\mathbf{B}_{-2}(X_{ij}) = (B_3(X_{ij}), \dots, B_L(X_{ij}))^T$. For a given λ , the AWG-Lasso objective function is

$$Q_A(\boldsymbol{\gamma}_{-1}; \lambda) = \frac{1}{n} \sum_{i=1}^n \rho_\tau(Y_i - \mu - \mathbf{W}_i^T \boldsymbol{\gamma}_{-1}) + \lambda \sum_{j=1}^p (w_{2j}|\boldsymbol{\gamma}_{2j}| + w_{-2j}|\boldsymbol{\gamma}_{-2j}|), \quad (\text{S.3})$$

where $\{(w_{2j}, w_{-2j})\}_{j=1}^p$ are obtained from some initial estimates. Minimizing $Q_A(\boldsymbol{\gamma}_{-1}; \lambda)$ w.r.t. $\boldsymbol{\gamma}_{-1}$, one gets

$$\hat{\boldsymbol{\gamma}}_{-1}^\lambda = \underset{\boldsymbol{\gamma}_{-1} \in R^{p(L-1)}}{\operatorname{argmin}} Q_A(\boldsymbol{\gamma}_{-1}; \lambda),$$

where $\hat{\boldsymbol{\gamma}}_{-1}^\lambda = (\hat{\gamma}_{21}^\lambda, \hat{\gamma}_{-21}^{\lambda T}, \dots, \hat{\gamma}_{2p}^\lambda, \hat{\gamma}_{-2p}^{\lambda T})^T$. Then, the model selected by AWG-Lasso is $\hat{\mathcal{S}}^\lambda = (\hat{\mathcal{S}}_l^\lambda, \hat{\mathcal{S}}_a^\lambda)$, where $\hat{\mathcal{S}}_l^\lambda = \{j \mid \hat{\gamma}_{2j}^\lambda \neq 0\}$ and $\hat{\mathcal{S}}_a^\lambda = \{j \mid \hat{\gamma}_{-2j}^\lambda \neq \mathbf{0}\}$. Like Section 2, this section also considers using HDIC to choose a suitable λ from a prescribed set Λ of positive numbers. Denote \mathbf{W}_i in (S.2) by $(v_{21i}, \mathbf{v}_{-21i}^T, \dots, v_{2pi}, \mathbf{v}_{-2pi}^T)^T$, where $(v_{2ji}, \mathbf{v}_{-2ji}^T)^T$ is the regressor vector corresponds to $\boldsymbol{\gamma}_{-1j}$. For a given model $\mathcal{S} = (\mathcal{S}_l, \mathcal{S}_a)$, define

$$R_A(\boldsymbol{\gamma}_\mathcal{S}) = \frac{1}{n} \sum_{i=1}^n \rho_\tau(Y_i - \mu - \mathbf{W}_{i\mathcal{S}}^T \boldsymbol{\gamma}_\mathcal{S}) \quad \text{and} \quad \tilde{\boldsymbol{\gamma}}_\mathcal{S} = \underset{\boldsymbol{\gamma}_\mathcal{S} \in R^{|\mathcal{S}_l| + (L-2)|\mathcal{S}_a|}}{\operatorname{argmin}} R_A(\boldsymbol{\gamma}_\mathcal{S}), \quad (\text{S.4})$$

where $\mathbf{W}_{i\mathcal{S}} \in R^{|\mathcal{S}_l| + (L-2)|\mathcal{S}_a|}$ consists of $\{v_{2ji} \mid j \in \mathcal{S}_l\}$ and $\{\mathbf{v}_{-2ji} \mid j \in \mathcal{S}_a\}$ and the corresponding coefficient $\boldsymbol{\gamma}_\mathcal{S} \in R^{|\mathcal{S}_l| + (L-2)|\mathcal{S}_a|}$ is conformably defined as in (6).

The HDIC value for model \mathcal{S} is stipulated by

$$\text{HDIC}(\mathcal{S}) = \log R_A(\tilde{\boldsymbol{\gamma}}_\mathcal{S}) + (|\mathcal{S}_l| + (L-2)|\mathcal{S}_a|) \frac{q_n \log p_n}{2n}, \quad (\text{S.5})$$

where p_n and q_n are defined as in Section 2. Let M_l and M_a be some known upper bounds for $|\mathcal{S}_l^0|$ and $|\mathcal{S}_a^0|$, respectively. We suggest choosing model $\hat{\mathcal{S}}^\lambda$, where

$$\hat{\lambda} = \underset{\lambda \in \Lambda, |\hat{\mathcal{S}}_l^\lambda| \leq M_l, |\hat{\mathcal{S}}_a^\lambda| \leq M_a}{\operatorname{argmin}} \text{HDIC}(\hat{\mathcal{S}}^\lambda).$$

S.2 Technical results for Theorems

S.2.1 Technical results for Theorem 1

We provide the proofs of Proposition 1 and Lemmas 1 and 2 here. We omit λ of $\hat{\boldsymbol{\gamma}}_\mathcal{S}^\lambda$ for notational simplicity.

First we state Lemma 3 for Proposition 1 and the notation for the lemma. Then we prove Proposition 1 by following Lemma 1 and Theorem 1 in Fan et al. (2014). Next we present the proofs of Lemmas 3, 1, and 2.

Before we state Lemma 3, we define

$$G_\mathcal{S}(M) = \sup_{\boldsymbol{\gamma}_\mathcal{S} \in \Gamma_\mathcal{S}(M)} |\{R_V(\boldsymbol{\gamma}_\mathcal{S}) - R_V(\boldsymbol{\gamma}_\mathcal{S}^*)\} - \mathbb{E}_\epsilon \{R_V(\boldsymbol{\gamma}_\mathcal{S}) - R_V(\boldsymbol{\gamma}_\mathcal{S}^*)\}|$$

where $\Gamma_\mathcal{S}(M)$ is defined in (31).

Lemma 3 Assume that Assumption A3 holds. For any fixed M , t , and \mathcal{S} , we have

$$P_\epsilon \left(G_{\mathcal{S}}(M) \geq 4M \sqrt{\frac{\Theta_1(\mathcal{S})}{n}} + t \right) \leq \exp \left\{ -\frac{nt^2}{8\Theta_1(\mathcal{S})M^2} \right\}.$$

When $t = K_0 M \{n^{-1}\Theta_1(\mathcal{S}) \log p_n\}^{1/2}$, we have from Lemma 3 that

$$P_\epsilon \left(G_{\mathcal{S}}(M) \geq (4 + K_0)M \sqrt{\frac{\Theta_1(\mathcal{S}) \log p_n}{n}} \right) \leq \exp(-K_0^2 \log p_n/8).$$

A remark is in place: A lower limit of the probability, $1 - \exp(-K_0^2 \log p_n/8)$, appears in the proof. But it is only related to evaluating $G_{\mathcal{S}}(M)$ and this $G_{\mathcal{S}}(M)$ does not contain the weights. Hence we can also deal with stochastic weights by using this proposition.

Proof of Proposition 1) We follow that of Theorem 1 in Fan et al. (2014). The following arguments do not depend on \mathcal{S} .

Taking $M = C_M LK_n(\mathcal{S})$, we evaluate the following expression on $\Gamma_{\mathcal{S}}(M)$.

$$E_\epsilon \{R_V(\gamma_{\mathcal{S}}) - R_V(\gamma_{\mathcal{S}}^*)\} = E_\epsilon \left[\frac{1}{n} \sum_{i=1}^n \{\rho_\tau(\epsilon'_i - a_i) - \rho_\tau(\epsilon'_i)\} \right], \quad (\text{S.6})$$

where we use the notation defined in (22) after Assumption A3 such as $\epsilon'_i = \epsilon_i + \delta_i$ and $a_i = \mathbf{W}_{i\mathcal{S}}^T(\gamma_{\mathcal{S}} - \gamma_{\mathcal{S}}^*)$. Note that

$$|a_i| \leq |\mathbf{W}_{i\mathcal{S}}|M \leq \Theta_5^{1/2}(\mathcal{S})M \rightarrow 0$$

due to the assumption of this proposition.

If $a_i > 0$, we have from the definition of $\rho_\tau(\cdot)$ that

$$\rho_\tau(\epsilon'_i - a_i) - \rho_\tau(\epsilon'_i) = \int_0^{a_i} I\{0 < \epsilon'_i \leq s\} ds + a_i(I\{\epsilon'_i \leq 0\} - \tau).$$

Then from Assumption A5, we obtain

$$\begin{aligned} & E_\epsilon \left[\int_0^{a_i} I\{0 < \epsilon'_i \leq s\} ds + a_i(I\{\epsilon'_i \leq 0\} - \tau) \right] \\ &= \int_0^{a_i} (F_i(s - \delta_i) - F_i(-\delta_i)) ds + a_i(\tau_i - \tau) \\ &= \frac{1}{2} f_i(-\delta_i) a_i^2 + o(a_i^2) + O(a_i^2 (\log n)^{-1}) + O(|\tau - \tau_i|^2 \log n). \end{aligned}$$

uniformly in i . Note that $|\tau - \tau_i|^2 \leq C_1 |\delta_i|^2$ for some positive C_1 and that we can deal with the case of $a_i < 0$ in the same way.

Hence the expression in (S.6) can be represented as

$$\frac{1}{2n} \sum_{i=1}^n f_i(-\delta_i) a_i^2 + o\left(n^{-1} \sum_{i=1}^n a_i^2\right) + O\left(n^{-1} \log n \sum_{i=1}^n \delta_i^2\right). \quad (\text{S.7})$$

The first term of (S.7) is written as

$$\begin{aligned} \frac{1}{2n} \sum_{i=1}^n f_i(-\delta_i) a_i^2 &= \frac{1}{2} (\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*)^T \frac{1}{n} \sum_{i=1}^n f_i(-\delta_i) \mathbf{W}_{iS} \mathbf{W}_{iS}^T (\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*) \\ &\geq \frac{\Theta_2(\mathcal{S})}{2L} |\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*|^2. \end{aligned} \quad (\text{S.8})$$

As for the third term of (S.7), we have from Assumption A3 that

$$\frac{\log n}{n} \sum_{i=1}^n \delta_i^2 = \frac{\log n}{n} \sum_{i=1}^n \left(\sum_{j \in \mathcal{S}_v^0} X_{ij} \delta_{ij} \right)^2 \leq \frac{\log n}{n} \sum_{i=1}^n \left(\sum_{j \in \mathcal{S}_v^0} X_{ij}^2 \right) \left(\sum_{j \in \mathcal{S}_v^0} \delta_{ij}^2 \right) \quad (\text{S.9})$$

$$\leq \frac{C_1 \log n}{nL^4} \sum_{i=1}^n \sum_{j \in \mathcal{S}_v^0} X_{ij}^2 \leq \frac{C_1 \log n}{L^4} \Theta_4 \quad (\text{S.10})$$

for some positive C_1 . We defined Θ_4 just before Assumption B4.

By combining (S.7), (S.8), and (S.9), we have

$$\mathbb{E}_\epsilon \{R_V(\boldsymbol{\gamma}_S) - R_V(\boldsymbol{\gamma}_S^*)\} \geq \frac{\Theta_2(\mathcal{S})}{2L} (1 + o(1)) |\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*|^2 + O\left(\frac{\Theta_4 \log n}{L^4}\right). \quad (\text{S.11})$$

We define $\boldsymbol{\gamma}_S^\alpha$ by

$$\boldsymbol{\gamma}_S^\alpha = \alpha \widehat{\boldsymbol{\gamma}}_S + (1 - \alpha) \boldsymbol{\gamma}_S^* \quad (\text{S.12})$$

for

$$0 \leq \alpha = \frac{M}{M + |\widehat{\boldsymbol{\gamma}}_S - \boldsymbol{\gamma}_S^*|} \leq 1.$$

Then

$$\boldsymbol{\gamma}_S^\alpha \in \Gamma_S(M).$$

Since the convexity of $Q_V(\boldsymbol{\gamma}_S)$ implies that

$$Q_V(\boldsymbol{\gamma}_S^\alpha) \leq \alpha Q_V(\widehat{\boldsymbol{\gamma}}_S) + (1 - \alpha) Q_V(\boldsymbol{\gamma}_S^*) \leq Q_V(\boldsymbol{\gamma}_S^*),$$

we have with probability larger than or equal to $1 - \exp(-K_0^2 \log p_n/8)$ that

$$\begin{aligned}
& \mathbb{E}_\epsilon[R_V(\gamma_S) - R_V(\gamma_S^*)]_{\gamma_S = \gamma_S^\alpha} \tag{S.13} \\
& \leq \frac{1}{n} \sum_{i=1}^n \rho_\tau(\gamma_S^*) - \mathbb{E}_\epsilon \left\{ \frac{1}{n} \sum_{i=1}^n \rho_\tau(\gamma_S^*) \right\} - \frac{1}{n} \sum_{i=1}^n \rho_\tau(\gamma_S^\alpha) + \mathbb{E}_\epsilon \left[\frac{1}{n} \sum_{i=1}^n \rho_\tau(\gamma_S) \right]_{\gamma_S = \gamma_S^\alpha} \\
& \quad + Q_V(\gamma_S^\alpha) - Q_V(\gamma_S^*) \\
& \quad - \lambda \sum_{j \in \mathcal{S}_c} w_{1j} |\gamma_{1j}^\alpha| - \lambda \sum_{j \in \mathcal{S}_v} w_{-1j} |\gamma_{-1j}^\alpha| + \lambda \sum_{j \in \mathcal{S}_c} w_{1j} |\gamma_{1j}^*| + \lambda \sum_{j \in \mathcal{S}_v} w_{-1j} |\gamma_{-1j}^*| \\
& \leq G_S(M) + \lambda |w_S| |\gamma_S^\alpha - \gamma_S^*| \\
& \leq (4 + K_0)M \left\{ \sqrt{\frac{\Theta_1(\mathcal{S}) \log p_n}{n}} + \lambda |w_S| \right\} = (4 + K_0)MK_n(\mathcal{S}).
\end{aligned}$$

By (S.11) and (S.13), we have

$$\begin{aligned}
|\gamma_S^\alpha - \gamma_S^*|^2 & \leq \frac{2(4 + K_0)L}{\Theta_2(\mathcal{S})} \{MK_n(\mathcal{S}) + O(\Theta_4 L^{-4} \log n)\} \\
& \leq \frac{2(4 + K_0)L}{\Theta_2(\mathcal{S})} \{C_M K_n^2(\mathcal{S})L + O(\Theta_4 L^{-4} \log n)\}
\end{aligned}$$

with probability larger than or equal to $1 - \exp(-K_0^2 \log p_n/8)$. Hence

$$\begin{aligned}
|\gamma_S^\alpha - \gamma_S^*| & \leq \frac{\{2(4 + K_0)\}^{1/2}}{\Theta_2^{1/2}(\mathcal{S})} \{C_M^{1/2} K_n(\mathcal{S})L + O(\Theta_4^{1/2} L^{-3/2} (\log n)^{1/2})\} \tag{S.14} \\
& \leq \frac{1}{2} C_M L K_n(\mathcal{S}) = \frac{1}{2} M
\end{aligned}$$

with probability larger than or equal to $1 - \exp(-K_0^2 \log p_n/8)$.

(S.12), (S.14), and simple algebra yield

$$|\widehat{\gamma}_S - \gamma_S^*| \leq M = C_M L K_n(\mathcal{S})$$

with probability larger than or equal to $1 - \exp(-K_0^2 \log p_n/8)$.

Hence the proof of the proposition is complete.

Proof of Lemma 3) We follow that of Lemma 1 in Fan et al. (2014).

Due to the Lipschitz continuity of $\rho_\tau(u)$ and application of the concentration inequalities (Theorems 14.3 and 14.4 in Bühlmann and van de Geer (2011)), we have

$$\begin{aligned}
\mathbb{E}_\epsilon \{G_S(M)\} & \leq 2\mathbb{E}_\epsilon \left[\sup_{\gamma_S \in \Gamma_S(M)} \left| \frac{1}{n} \sum_{i=1}^n \xi_i \{ \rho_\tau(Y_i - \mathbf{W}_{iS}^T \gamma_S) - \rho_\tau(Y_i - \mathbf{W}_{iS}^T \gamma_S^*) \} \right| \right] \\
& \leq 4\mathbb{E}_\epsilon \left[\sup_{\gamma_S \in \Gamma_S(M)} \left| \frac{1}{n} \sum_{i=1}^n \xi_i \mathbf{W}_{iS}^T (\gamma_S - \gamma_S^*) \right| \right],
\end{aligned}$$

where $\{\xi_j\}_{j=1}^n$ is a Rademacher sequence of and independent of $\{(Y_j, \mathbf{X}_j, Z_j)\}_{j=1}^n$. Since

$$\begin{aligned}
& \left| \sum_{i=1}^n \xi_i \mathbf{W}_{iS}^T (\gamma_S - \gamma_S^*) \right| \\
&= \left| \sum_{j \in \mathcal{S}_c} \left(\sum_{i=1}^n \xi_i X_{ij} L^{-1/2} \right) (\gamma_{1j} - \gamma_{1j}^*) + \sum_{j \in \mathcal{S}_v} \left\{ \sum_{i=1}^n \xi_i X_{ij} \mathbf{B}_{-1}^T(Z_i) (\gamma_{-1j} - \gamma_{-1j}^*) \right\} \right| \\
&\leq |\gamma_S - \gamma_S^*| \left\{ \sum_{j \in \mathcal{S}_c} \left| \sum_{i=1}^n \xi_i X_{ij} L^{-1/2} \right|^2 + \sum_{j \in \mathcal{S}_v} \left| \sum_{i=1}^n \xi_i X_{ij} \mathbf{B}_{-1}(Z_i) \right|^2 \right\}^{1/2},
\end{aligned}$$

we have

$$\begin{aligned}
& \mathbb{E}_\epsilon \{G_S(M)\} \tag{S.15} \\
&\leq \frac{4M}{n^{1/2}} \mathbb{E}_\epsilon \left[\left\{ \frac{1}{n} \sum_{j \in \mathcal{S}_c} \left| \sum_{i=1}^n \xi_i X_{ij} L^{-1/2} \right|^2 + \frac{1}{n} \sum_{j \in \mathcal{S}_v} \left| \sum_{i=1}^n \xi_i X_{ij} \mathbf{B}_{-1}(Z_i) \right|^2 \right\}^{1/2} \right] \\
&\leq \frac{4M}{n^{1/2}} \left[\mathbb{E}_\epsilon \left\{ \frac{1}{n} \sum_{j \in \mathcal{S}_c} \left| \sum_{i=1}^n \xi_i X_{ij} L^{-1/2} \right|^2 + \frac{1}{n} \sum_{j \in \mathcal{S}_v} \left| \sum_{i=1}^n \xi_i X_{ij} \mathbf{B}_{-1}(Z_i) \right|^2 \right\} \right]^{1/2} \\
&\leq \frac{4M}{n^{1/2}} \left\{ \frac{1}{n} \sum_{i=1}^n |\mathbf{W}_{iS}|^2 \right\}^{1/2} \leq 4M \sqrt{\frac{\Theta_1(\mathcal{S})}{n}}.
\end{aligned}$$

Next we apply Massart's inequality (Theorem 14.2 in Bühlmann and van de Geer (2011)) to evaluate the stochastic part $G_S(M) - \mathbb{E}_\epsilon \{G_S(M)\}$. Then noticing

$$|\mathbf{W}_{iS}^T (\gamma_S - \gamma_S^*)|^2 \leq |\mathbf{W}_{iS}|^2 |\gamma_S - \gamma_S^*|^2 \leq |\mathbf{W}_{iS}|^2 M^2$$

and

$$\frac{1}{n} \sum_{i=1}^n |\mathbf{W}_{iS}|^2 M^2 \leq \Theta_1(\mathcal{S}) M^2,$$

we have as in Lemma 1 in Fan et al. (2014)

$$\mathbb{P}_\epsilon \left(G_S(M) \geq 4M \sqrt{\frac{\Theta_1(\mathcal{S})}{n}} + t \right) \leq \exp \left\{ - \frac{nt^2}{8\Theta_1(\mathcal{S})M^2} \right\}.$$

We used (S.15) to evaluate $\mathbb{E}_\epsilon \{G_S(M)\}$ in the conditional probability.

Hence the proof of the lemma is complete.

Proof of Lemma 1) Recall that $\mathbf{B}(z) = A_0 \mathbf{B}_0(z)$ and note (S.48) in Section S.3. Thus we have only to demonstrate

$$\left| \frac{1}{n} \sum_{i=1}^n B_{0l}(Z_i) X_{ij} \rho'_\tau(\epsilon_i + \delta_i) \right| \leq C_1 \{(nL)^{-1} \log p_n\}^{1/2} \tag{S.16}$$

uniformly in l and j with probability tending to 1 for some positive C_1 . Recall $B_{0l}(z)$ is the l -th element of the B-spline basis.

Note that

$$\mathbb{E}_\epsilon \left\{ \frac{1}{n} \sum_{i=1}^n B_{0l}(Z_i) X_{ij} \rho'_\tau(\epsilon_i + \delta_i) \right\} = \frac{1}{n} \sum_{i=1}^n B_{0l} X_{ij}(Z_i) (\tau - \tau_i)$$

and $|\tau - \tau_i| = O(L^{-2})$ uniformly in i .

Since Assumption A4 implies

$$\mathbb{E} \left\{ \frac{1}{n} \sum_{i=1}^n B_{0l}(Z_i) X_{ij} (\tau - \tau_i) \right\} = O(L^{-3})$$

and

$$\text{Var} \left\{ \frac{1}{n} \sum_{i=1}^n B_{0l}(Z_i) X_{ij} (\tau - \tau_i) \right\} = O(n^{-1} L^{-5}),$$

uniformly in l and j , we apply Bernstein's inequality unconditionally and obtain

$$\left| \mathbb{E}_\epsilon \left\{ \frac{1}{n} \sum_{i=1}^n B_{0l}(Z_i) X_{ij} \rho'_\tau(\epsilon_i + \delta_i) \right\} \right| \leq C_2 \{(nL^5)^{-1} \log p_n\}^{1/2} + O(L^{-3}) \quad (\text{S.17})$$

uniformly in l and j with probability tending to 1 for some positive C_2 .

Noticing that

$$\frac{1}{n} \sum_{i=1}^n B_{0l}^2(Z_i) X_{ij}^2 \leq C_3 L^{-1}$$

uniformly in l and j with probability tending to 1 for some positive C_3 , we apply Bernstein's inequality conditionally and obtain

$$\left| \frac{1}{n} \sum_{i=1}^n B_{0l}(Z_i) X_{ij} \rho'_\tau(\epsilon_i + \delta_i) - \mathbb{E}_\epsilon \left\{ \frac{1}{n} \sum_{i=1}^n B_{0l}(Z_i) X_{ij} \rho'_\tau(\epsilon_i + \delta_i) \right\} \right| \leq C_4 \{(nL)^{-1} \log p_n\}^{1/2} \quad (\text{S.18})$$

uniformly in l and j with probability tending to 1 for some positive C_4 .

Hence (S.16) follows from (S.17) and (S.18) and the proof of the lemma is complete.

Proof of Lemma 2) We can prove this lemma almost in the same way as Lemma B.5 in Sherwood and Wang (2016) and the detailed proof is very lengthy. We just outline the proof.

First we define $d_{lj}(\gamma_{S^0})$ by

$$\begin{aligned} d_{lj}(\gamma_{S^0}) = & \frac{1}{n} \sum_{i=1}^n B_{0l}(Z_i) X_{ij} [\rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}) - \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}^*) \\ & - \mathbb{E}_\epsilon \{ \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}) - \rho'_\tau(Y_i - \mathbf{W}_{iS^0}^T \gamma_{S^0}^*) \}] \end{aligned}$$

and take and fix any positive C_0 . Then as in the proof of Lemma 1, we have only to prove that

$$|d_{lj}(\gamma_{S^0})| \leq C_1 \{(nL \log n)^{-1} \log p_n\}^{1/2}$$

uniformly in $l, j \in \overline{\mathcal{S}_v^0}$, and $\gamma_{S^0} \in \Gamma_{S^0}(C_0 L(n^{-1} \log p_n)^{1/2} (\log n)^k)$ with probability tending to 1 for some positive C_1 depending on C_0 .

Note that the conditional variance of $d_{lj}(\gamma_{S^0})$ is uniformly bounded by

$$\frac{C_2 X_M}{nL} L(n^{-1} \log p_n)^{1/2} (\log n)^k \leq C_3 X_M \{n^{-3} (\log n)^{2k} \log p_n\}^{1/2}$$

with probability tending to 1 for some positive C_2 and C_3 . They depend on C_0 . Besides, we can cover $\Gamma_{S^0}(C_0 L(n^{-1} \log p_n)^{1/2} (\log n)^k)$ by N open balls with radius

$$[\{C_0 L(n^{-1} \log p_n)^{1/2} (\log n)^k\} n^{-2m}]^{1/2}$$

for any large fixed m and this N satisfies

$$N = O(n^{md_V(\mathcal{S}^0)}).$$

See Lemma 2.5 in van de Geer (2000) for this upper bound of N . We denote the centers of the covering open balls by $\gamma_1, \dots, \gamma_N$. Note that

$$pLN = O(\exp\{\log p_n + md_V(\mathcal{S}^0) \log n\}).$$

For any γ_s among the centers, we have by employing Bernstein's inequality conditionally that

$$P_\epsilon \left(|d_{lj}(\gamma_s)| \geq C_4 \sqrt{\frac{\log p_n}{nL \log n}} \right) \leq \exp \left\{ -C_3 \frac{(\log p_n)^{1/2} n^{3/10}}{X_M (\log n)^{k+1}} \right\}$$

uniformly in γ_s with probability tending to 1 for some positive C_4 and C_5 and we also have from Assumption A4 that

$$pLN \exp \left\{ -C_3 \frac{(\log p_n)^{1/2} n^{3/10}}{X_M (\log n)^{k+1}} \right\} = \exp \left[C_6 \{\log p_n + md_V(\mathcal{S}^0) \log n\} - C_3 \frac{(\log p_n)^{1/2} n^{3/10}}{X_M (\log n)^{k+1}} \right] \\ \rightarrow 0$$

for some positive C_6 . Therefore we successfully evaluated $d_{lj}(\gamma_{S^0})$ at all the centers.

We can evaluate $d_{lj}(\gamma_{S^0})$ inside the open balls exactly as in the proof of Lemma B.5 in Sherwood and Wang (2016) since we can take any large m . Hence the proof of the lemma is complete.

S.2.2 Technical results for Theorem 2

In this subsection, we state Lemma 4 and then give the proofs of the underfitting case, (44), and Lemma 4.

First we state Lemma 4, which is used to evaluate the bias from $(\tau_i - \tau)$ in the proof of Theorem 2. Note that the Hölder continuity of g_j'' with exponent α is almost sufficient for $\tau_i - \tau = O_p(X_M L^{-(2+\alpha)})$.

Recall the definition of \mathbf{b}_S in (45) and \mathbf{b}_{S_2} is defined as \mathbf{a}_{S_2} in (52). By using the properties of \mathbf{b}_S and \mathbf{b}_{S_2} in this lemma and replacing \mathbf{a}_S with $\mathbf{a}_S + \mathbf{b}_S$ in (45), we can prove Theorem 2 in the same way if $X_M^4 L^{-2\alpha} = O(L^{-1})$. Recall that $L = c_L n^{1/5}$ in this paper. Both of $|\mathbf{b}_S|^2$ and $|\mathbf{b}_{S_2}|^2$ have $O_p\left(\frac{X_M^4 \log n}{L^{5+2\alpha}}\right)$ and these are not typos.

Lemma 4 *In the setup of Theorem 2, we have*

$$(\gamma_S - \gamma_S^*)^T \mathbf{b}_S = O_p\left(\frac{(q_n \log p_n)^{1/2}}{n}\right) \quad (\text{S.19})$$

uniformly in $\gamma_S \in \Gamma_S(M_1 L (q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} for any fixed M_1 . Let Assumption A3' be replaced with Assumption A3. If $\tau_i - \tau = O_p(X_M L^{-(2+\alpha)})$ uniformly in i for some nonnegative α , we have

$$|\mathbf{b}_S|^2 = O_p\left(\frac{X_M^4 \log n}{L^{5+2\alpha}}\right) \quad \text{and} \quad |\mathbf{b}_{S_2}|^2 = O_p\left(\frac{X_M^4 \log n}{L^{5+2\alpha}}\right)$$

uniformly in \mathcal{S} .

Since $d_V(\mathcal{S}) = |\mathcal{S}_c| + (L-1)|\mathcal{S}_v|$, we have uniformly in \mathcal{S} ,

$$|\mathbf{b}_S|^2 = O_p\left(\frac{\log p_n}{n}\right) \quad \text{and} \quad |\mathbf{b}_{S_2}|^2 = (d_V(\mathcal{S}) - d_V(\mathcal{S}^0)) O_p\left(\frac{\log p_n}{nL}\right)$$

as in (47) and (55) if $X_M^4 L^{-2\alpha} = O(L^{-1})$. Then we can prove Theorem 2 in the same way.

Proof of the underfitting case) Next we consider the underfitting case. For $\mathcal{S} = (\mathcal{S}_c, \mathcal{S}_v)$ that does not include \mathcal{S}^0 and satisfies

$$|\mathcal{S}_c| \leq M_c \quad \text{and} \quad |\mathcal{S}_v| \leq M_v,$$

we put

$$\mathcal{S}^+ = \mathcal{S} \cup \mathcal{S}^0. \quad (\text{S.20})$$

Then $\mathcal{S}^+ \in \mathcal{S}_{C+M}$ in (26). Note that uniform results proved in the overfitting case still hold for \mathcal{S}^+ in (S.20).

Since

$$\log R_V(\tilde{\gamma}_S) - \log R_V(\tilde{\gamma}_{S^0}) = \log \left\{ 1 + \frac{R_V(\tilde{\gamma}_S) - R_V(\tilde{\gamma}_{S^0})}{R_V(\tilde{\gamma}_{S^0})} \right\}$$

and

$$R_V(\tilde{\gamma}_{S^0}) = \frac{1}{n} \sum_{i=1}^n \rho_\tau(\epsilon_i) + o_p(1) = \mathbb{E}\{\rho_\tau(\epsilon_i)\} + o_p(1), \quad (\text{S.21})$$

we have only to demonstrate

$$R_V(\tilde{\gamma}_S) - R_V(\tilde{\gamma}_{S^0}) > C_1 L \zeta_n^2 \frac{\log p_n}{2n} \quad (\text{S.22})$$

uniformly in \mathcal{S} with probability tending to 1 for some C_1 and ζ_n such that $\zeta_n/q_n^{1/2} = C_\zeta$. Note that we should be able to take and fix any sufficiently large C_ζ and that C_1 has to be independent of C_ζ when C_ζ is large. Then Assumption A1 and (S.21) assure (S.22) dominates the penalty terms. Since (S.21) follows from the argument for the overfitting case and Assumption A5', we consider only (S.22).

From Assumption A2', we have uniformly in \mathcal{S} ,

$$\frac{|\gamma_{S^0-\mathcal{S}}^*|}{L(n^{-1}q_n \log p_n)^{1/2}} \rightarrow \infty,$$

where $\gamma_{S^0-\mathcal{S}}^*$ is obtained by removing all the j -th elements satisfying $j \in \mathcal{S} \cap \mathcal{S}^0$ from $\gamma_{S^0}^*$.

Since \mathcal{S}^+ includes \mathcal{S}^0 and \mathcal{S} does not include \mathcal{S}^0 , Proposition 1 with no penalty implies that

$$|(\tilde{\gamma}_S^T, \mathbf{0}^T)^T - \tilde{\gamma}_{S^+}| > L \zeta_n (n^{-1} \log p_n)^{1/2} \quad (\text{S.23})$$

uniformly in \mathcal{S} with probability tending to 1 for $\zeta_n = C_\zeta q_n^{1/2}$. Note that we can take and fix any large C_ζ here. This also holds with $\tilde{\gamma}_{S^+}$ replaced by $\bar{\gamma}_{S^+}$ due to (51).

Let us follow the standard arguments for general underfitting cases. There is an $0 < \alpha < 1$ such that

$$|\alpha((\tilde{\gamma}_S^T, \mathbf{0}^T)^T - \bar{\gamma}_{S^+})| = L \zeta_n (n^{-1} \log p_n)^{1/2}$$

and set

$$\Delta_S = \alpha((\tilde{\gamma}_S^T, \mathbf{0}^T)^T - \bar{\gamma}_{S^+}).$$

The arguments from (44) to (50) imply that

$$\begin{aligned}
R_V(\bar{\gamma}_{\mathcal{S}^+} + \Delta_{\mathcal{S}}) &\geq R_V(\bar{\gamma}_{\mathcal{S}^+}) + C_2 \zeta_n^2 \frac{L \log p_n}{2n} + O_p\left(\frac{(q_n \log p_n)^{1/2}}{n}\right) + O_p\left(\frac{\log p_n}{n(\log n)^2}\right) \\
&\geq R_V(\bar{\gamma}_{\mathcal{S}^+}) + C_2 \zeta_n^2 \frac{L \log p_n}{4n} \geq R_V(\tilde{\gamma}_{\mathcal{S}^+}) + C_2 \zeta_n^2 \frac{L \log p_n}{4n}
\end{aligned} \tag{S.24}$$

uniformly in \mathcal{S} with probability tending to 1 for some positive C_2 independent of C_ζ . We used the optimality of $\tilde{\gamma}_{\mathcal{S}^+}$ and Assumption B5 here.

Because of (S.24), the convexity of $R_V(\gamma_{\mathcal{S}^+})$, and the definition of $\Delta_{\mathcal{S}}$, we have

$$R_V(\tilde{\gamma}_{\mathcal{S}}) \geq R_V(\bar{\gamma}_{\mathcal{S}^+} + \Delta_{\mathcal{S}}) \geq R_V(\bar{\gamma}_{\mathcal{S}^+}) \geq R_V(\tilde{\gamma}_{\mathcal{S}^+}) \tag{S.25}$$

uniformly in \mathcal{S} with probability tending to 1. From (S.24) and (S.25), we obtain

$$R_V(\tilde{\gamma}_{\mathcal{S}}) \geq R_V(\tilde{\gamma}_{\mathcal{S}^+}) + C_2 \zeta_n^2 \frac{L \log p_n}{4n} \tag{S.26}$$

uniformly in \mathcal{S} with probability tending to 1. Recalling the results for the overfitting case such as (50) and the evaluation of (53), we have

$$R_V(\tilde{\gamma}_{\mathcal{S}^+}) \geq R_V(\tilde{\gamma}_{\mathcal{S}^0}) + (d_V(\mathcal{S}^0) - d(\mathcal{S}^+)) \frac{q_n \log p_n}{2n} \tag{S.27}$$

uniformly in \mathcal{S} with probability tending to 1.

By combining (S.26) and (S.27), we get

$$R_V(\tilde{\gamma}_{\mathcal{S}}) \geq R_V(\tilde{\gamma}_{\mathcal{S}^0}) + C_2 \zeta_n^2 \frac{L \log p_n}{4n} + (d_V(\mathcal{S}^0) - d(\mathcal{S}^+)) \frac{q_n \log p_n}{2n} \tag{S.28}$$

uniformly in \mathcal{S} with probability tending to 1. Since $d_V(\mathcal{S}^0) - d(\mathcal{S}^+) = O(L)$ from Assumption A1 and $\zeta_n = C_\zeta q_n^{1/2}$, we have from (S.28) that

$$R_V(\tilde{\gamma}_{\mathcal{S}}) > R_V(\tilde{\gamma}_{\mathcal{S}^0}) + C_3 \zeta_n^2 \frac{L \log p_n}{2n} \tag{S.29}$$

for any sufficiently large fixed C_ζ uniformly in \mathcal{S} with probability tending to 1. Note that C_3 is independent of C_ζ when C_ζ is larger than some value depending on the assumptions.

Hence we have established (S.22) and the proof of the underfitting case is complete.

Proof of (44) We take a positive M_1 and consider

$$\begin{aligned}
& R_V(\boldsymbol{\gamma}_S) - R_V(\boldsymbol{\gamma}_S^*) & (S.30) \\
& + \frac{1}{n} \sum_{i=1}^n \mathbf{W}_{iS}^T(\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*)(\tau - I\{\epsilon'_i \leq 0\}) - \mathbb{E}_\epsilon\{R_V(\boldsymbol{\gamma}_S) - R_V(\boldsymbol{\gamma}_S^*)\} \\
& - \frac{1}{n} \sum_{i=1}^n \mathbf{W}_{iS}^T(\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*)(\tau - \tau_i) \\
& = \frac{1}{n} \sum_{i=1}^n D_i(\boldsymbol{\gamma}_S),
\end{aligned}$$

where $D_i(\boldsymbol{\gamma}_S)$ is clearly defined in the above equation, $\tau_i = \mathbb{P}_\epsilon(\epsilon'_i \leq 0)$, and $|\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*| \leq M_1 L(q_n n^{-1} \log p_n)^{1/2}$.

We show that

$$\frac{1}{n} \sum_{i=1}^n D_i(\boldsymbol{\gamma}_S) = O_p\left(\frac{\log p_n}{n(\log n)^2}\right) \quad (S.31)$$

uniformly in $\boldsymbol{\gamma}_S \in \Gamma_S(M_1 L(q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} for any fixed M_1 . To verify (S.31), we should note that

$$D_i(\boldsymbol{\gamma}_S) = \bar{D}_i(\boldsymbol{\gamma}_S) - \mathbb{E}_\epsilon\{\bar{D}_i(\boldsymbol{\gamma}_S)\}, \quad (S.32)$$

where

$$\bar{D}_i(\boldsymbol{\gamma}_S) = \rho_\tau(Y_i - \mathbf{W}_{iS}^T \boldsymbol{\gamma}_S) - \rho_\tau(\mathbf{W}_{iS}^T \boldsymbol{\gamma}_S^*) + \mathbf{W}_{iS}^T(\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*)(\tau - I\{\epsilon'_i \leq 0\})$$

and that

$$\rho_\tau(\epsilon'_i - a_i) - \rho_\tau(\epsilon'_i) = -a_i(\tau - I\{\epsilon'_i \leq 0\}) - (\epsilon'_i - a_i)[I\{\epsilon'_i \leq a_i\} - I\{\epsilon'_i \leq 0\}], \quad (S.33)$$

where $a_i = \mathbf{W}_{iS}^T(\boldsymbol{\gamma}_S - \boldsymbol{\gamma}_S^*)$.

By using (S.33), we can obtain the following three facts (S.34)-(S.36) uniformly in $\boldsymbol{\gamma}_S \in \Gamma_S(M_1 L(q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} . Note that C_2, \dots, C_7 are some positive constants.

$$\max_{1 \leq i \leq n} |\mathbf{W}_{iS}| \leq C_2 X_M (M_c^{1/2} L^{-1/2} + M_v^{1/2}) \leq C_3 X_M \quad (\text{S.34})$$

$$\max_{1 \leq i \leq n} |\bar{D}_i(\gamma_S)| \leq \max_{1 \leq i \leq n} |\mathbf{W}_{iS}| M_1 L (q_n n^{-1} \log p_n)^{1/2} \leq C_4 X_M M_1 L (q_n n^{-1} \log p_n)^{1/2} \quad (\text{S.35})$$

$$\begin{aligned} \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}_\epsilon \{ [\bar{D}_i(\gamma_S)]^2 \} &\leq \frac{C_5}{n^2} \sum_{i=1}^n |\mathbf{W}_{iS}^T (\gamma_S - \gamma_S^*)|^3 \\ &\leq \frac{C_6}{n} \max_{1 \leq i \leq n} |\mathbf{W}_{iS}| \{M_1 L (q_n n^{-1} \log p_n)^{1/2}\}^3 \lambda_{\max} \left(n^{-1} \sum_{i=1}^n \mathbf{W}_{iS} \mathbf{W}_{iS}^T \right) \\ &\leq \frac{C_7 M_1^3 M_2}{n} L^2 X_M (q_n n^{-1} \log p_n)^{3/2} \end{aligned} \quad (\text{S.36})$$

if

$$\lambda_{\max} \left(n^{-1} \sum_{i=1}^n \mathbf{W}_{iS} \mathbf{W}_{iS}^T \right) \leq \frac{M_2}{L}. \quad (\text{S.37})$$

By using (S.34)-(S.36) and Bernstein's inequality, we have

$$\mathbb{P}_\epsilon \left(\left| n^{-1} \sum_{i=1}^n D_i(\gamma_S) \right| \geq \frac{\log p_n}{n(\log n)^2} \right) \leq C_8 \exp \left\{ - \frac{C_9 n^{1/10} (\log p_n)^{1/2}}{M_1^3 M_2 q_n^{3/2} X_M (\log n)^4} \right\} \quad (\text{S.38})$$

for any fixed $\gamma_S \in \Gamma_S(M_1 L (q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} if (S.37) holds. Note that C_8 and C_9 are some positive constants.

By appealing to the standard argument based on the Lipschitz continuity and (S.38) and using Assumptions A4' and B5, we obtain (S.31) uniformly in $\gamma_S \in \Gamma_S(M_1 L (q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} for any fixed M_1 .

We evaluate $\mathbb{E}_\epsilon \{ R_V(\gamma_S) - R_V(\gamma_S^*) \}$ in (S.30) by using (S.33) and Assumption A5'. Since

$$\mathbb{E}_\epsilon \{ \rho_\tau(\epsilon'_i - a_i) - \rho_\tau(\epsilon'_i) \} = \frac{1}{2} f_i(-\delta_i) a_i^2 + a_i(\tau_i - \tau) + O(|a_i|^3),$$

where $a_i = \mathbf{W}_{iS}^T (\gamma_S - \gamma_S^*)$, we have

$$\mathbb{E}_\epsilon \{ R_V(\gamma_S) - R_V(\gamma_S^*) \} = \frac{1}{2n} \sum_{i=1}^n f_i(-\delta_i) a_i^2 + \frac{1}{n} \sum_{i=1}^n a_i(\tau_i - \tau) + O\left(\frac{1}{n} \sum_{i=1}^n |a_i|^3 \right) \quad (\text{S.39})$$

uniformly in $\gamma_S \in \Gamma_S(M_1 L (q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} . Assumption A4' implies that uniformly in $\gamma_S \in \Gamma_S(M_1 L (q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} for any fixed M_1 ,

$$\frac{1}{n} \sum_{i=1}^n |a_i|^3 \leq \frac{\max_{i=1}^n |a_i|}{n} \sum_{i=1}^n |a_i|^2 = O_p\left(\frac{\log p_n}{n(\log n)^2} \right). \quad (\text{S.40})$$

By (S.39) and (S.40), we obtain

$$\begin{aligned} E_\epsilon\{R_V(\gamma_S) - R_V(\gamma_S^*)\} &= \frac{1}{2}(\gamma_S - \gamma_S^*)^T \widehat{\Sigma}_S(\gamma_S - \gamma_S^*) \\ &\quad + \frac{1}{n} \sum_{i=1}^n a_i(\tau_i - \tau) + O_p\left(\frac{\log p_n}{n(\log n)^2}\right) \end{aligned} \quad (\text{S.41})$$

uniformly in $\gamma_S \in \Gamma_S(M_1 L(q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} for any fixed M_1 .

By combining (S.30), (S.31), and (S.41), we obtain (44),

$$\begin{aligned} R_V(\gamma_S) - R_V(\gamma_S^*) &= -(\gamma_S - \gamma_S^*)^T \frac{1}{n} \sum_{i=1}^n \mathbf{W}_{iS}(\tau_i - I\{\epsilon'_i \leq 0\}) + \frac{1}{2}(\gamma_S - \gamma_S^*)^T \widehat{\Sigma}_S(\gamma_S - \gamma_S^*) \\ &\quad + (\gamma_S - \gamma_S^*)^T \frac{1}{n} \sum_{i=1}^n \mathbf{W}_{iS}(\tau_i - \tau) + O_p\left(\frac{\log p_n}{n(\log n)^2}\right) \end{aligned}$$

uniformly in $\gamma_S \in \Gamma_S(M_1 L(q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} for any fixed M_1 . Hence the proof of (44) is complete.

Proof of Lemma 4) We prove the former half by using Assumption A3'. By exploiting (23) and Assumptions A3' and B4', we have

$$\frac{1}{n} \sum_{i=1}^n |a_i(\tau_i - \tau)| \leq \left(n^{-1} \sum_{i=1}^n a_i^2\right)^{1/2} \left(n^{-1} \sum_{i=1}^n (\tau_i - \tau)^2\right)^{1/2} = O_p\left(\frac{(q_n \log p_n)^{1/2}}{n}\right).$$

uniformly in $\gamma_S \in \Gamma_S(M_1 L(q_n n^{-1} \log p_n)^{1/2})$ and \mathcal{S} since

$$\frac{1}{n} \sum_{i=1}^n a_i^2 = O_p(n^{-1} L q_n \log p_n) \quad \text{and} \quad \frac{1}{n} \sum_{i=1}^n (\tau_i - \tau)^2 = O_p(L^{-6})$$

uniformly as well.

Before we consider the latter, we should recall that $\mathbf{B}(z) = A_0 \mathbf{B}_0(z)$, where $\mathbf{B}_0(z) = (B_{01}(z), \dots, B_{0L}(z))^T$ is the equispaced B-spline basis on $[0, 1]$, and that the first element of $\mathbf{B}(z)$ is $L^{-1/2}$. Therefore we should deal with

$$\frac{X_M}{nL^{1/2}} \sum_{i=1}^n |\tau_i - \tau| \quad (\text{S.42})$$

and

$$\frac{X_M}{n} \sum_{i=1}^n \mathbf{B}_0(Z_i) |\tau_i - \tau|. \quad (\text{S.43})$$

As for (S.42), we have

$$\frac{X_M}{nL^{1/2}} \sum_{i=1}^n |\tau_i - \tau| = O_p\left(\frac{X_M^2}{L^{2+1/2+\alpha}}\right) \quad (\text{S.44})$$

from the assumption on $(\tau_i - \tau)$.

Since we have $E\{B_{0j}(Z_i)\} = O(L^{-1})$ uniformly in j , we have

$$\frac{X_M}{n} \sum_{i=1}^n B_{0j}(Z_i) |\tau_i - \tau| = O_p\left(\frac{X_M^2 (\log n)^{1/2}}{L^{3+\alpha}}\right) \quad (\text{S.45})$$

uniformly in j from the standard argument based on Bernstein's inequality.

(S.44) and (S.45) yields that

$$|\mathbf{b}_S|^2 = |\mathcal{S}_v| O_p\left(\frac{X_M^4 \log n}{L^{5+2\alpha}}\right) + |\mathcal{S}_c| O_p\left(\frac{X_M^4}{L^{5+2\alpha}}\right)$$

uniformly in \mathcal{S} .

The result for \mathbf{b}_{S_2} follows from the same argument. Hence the proof of the lemma is complete.

S.3 Properties of B-spline bases

We describe properties of our basis and give comments on some misleading assumptions on spline bases in the literature for reference.

First we describe how to construct our orthonormal spline basis $\mathbf{B}(z) = (B_1(z), \dots, B_L(z))^T$ from the equispaced B-spline basis on $[0, 1]$, which is denoted by $\mathbf{B}_0(z) = (B_{01}(z), \dots, B_{0L}(z))^T$. Recall that $L = c_L n^{1/5}$ in this paper. We also should recall two well-known facts:

$$\sum_{j=1}^L B_{0j}(z) = 1 \quad \text{and} \quad B_{0j}(z) \geq 0 \quad (\text{S.46})$$

$$\frac{C_1}{L} \leq \lambda_{\min}(\Omega_0) \leq \lambda_{\max}(\Omega_0) \leq \frac{C_2}{L} \quad (\text{S.47})$$

where $\Omega_0 = \int_0^1 \mathbf{B}_0(z) \mathbf{B}_0^T(z) dz$ and C_1 and C_2 are positive constants and independent of L .

Therefore there exists an $L \times L$ matrix A_0 such that

$$\begin{aligned} \mathbf{B}(z) &= A_0 \mathbf{B}_0(z), \quad \int_0^1 \mathbf{B}(z) \mathbf{B}^T(z) dz = A_0 \Omega_0 A_0^T = L^{-1} I_L, \\ B_1(z) &= L^{1/2}, \quad \text{and} \quad B_2(z) = \sqrt{\frac{12}{L}} \left(z - \frac{1}{2}\right). \end{aligned}$$

We denote the $L \times L$ identity matrix by I_L .

We can obtain an A_0 numerically by carrying out the Gram-Schmidt orthonormalization. Notice also that

$$C_3 \leq \lambda_{\min}(A_0 A_0^T) \leq \lambda_{\max}(A_0 A_0^T) \leq C_4, \quad (\text{S.48})$$

where C_3 and C_4 are positive constants and independent of L .

When we deal with varying coefficient models, $B_1(z) = L^{-1/2}$ is used for the constant parts and $\mathbf{B}_{-1} = (B_2(z), \dots, B_L(z))^T$ is used for the non-constant parts. When we deal with additive models, $B_2(z) = \sqrt{\frac{12}{L}}\left(z - \frac{1}{2}\right)$ is used for the linear parts and $(B_3(z), \dots, B_L(z))^T$ is used for the nonlinear parts.

Next we consider approximation by our spline basis $\mathbf{B}(z) = (B_1(z), \mathbf{B}_{-1}^T(z))^T = (B_1(z), B_2(z), \mathbf{B}_{-2}^T(z))^T$ under Assumption A3. Assume that

$$\|g\|_{\infty} + \|g'\|_{\infty} + \|g''\|_{\infty} \leq C_g.$$

Varying coefficient models: There exists $\gamma_{-1}^* \in R^{L-1}$ such that $\|g_n - \gamma_{-1}^{*T} \mathbf{B}_{-1}\|_{\infty} \leq C_1 C_g L^{-2}$. We can take $\gamma_1^* = L^{1/2} g_c$.

Additive models: Let $g(x)$ satisfy $\int_0^1 g(x) dx = 0$. Then there exist $\gamma_2^* \in R$ and $\gamma_{-2}^* \in R^{L-2}$ such that

$$\|g_l - \gamma_2^* B_2\|_{\infty} + \|g_a - \gamma_{-2}^{*T} \mathbf{B}_{-2}\|_{\infty} \leq C_2 C_g L^{-2}.$$

Note that C_1 and C_2 are independent of the specific function. We verify the latter here since the former is easier.

Corollary 6.26 in Schumaker (2007) implies that there is $\boldsymbol{\gamma}^* = (\gamma_1^*, \gamma_2^*, \gamma_{-2}^{*T})^T$ such that

$$\|g - \boldsymbol{\gamma}^{*T} \mathbf{B}\|_{\infty} \leq C_3 C_g L^{-2} \quad (\text{S.49})$$

since $\mathbf{B}(x)$ is constructed from $\mathbf{B}_0(x)$. Noticing

$$\gamma_1^* = L^{1/2} \int_0^1 (\boldsymbol{\gamma}^{*T} \mathbf{B}(x) - g(x)) dx$$

and $|\gamma_1^*| \leq C_3 C_g L^{-3/2}$, we can take $\gamma_1^* = 0$ without affecting (S.49).

Put

$$g^*(x) = \gamma_2^* B_2(x) + \gamma_{-2}^{*T} \mathbf{B}_{-2}(x) \quad \text{and} \quad g_l(x) = \gamma_2' B_2(x) + g_a(x),$$

where γ_2' is defined in the second equation and $g_l(x) = \gamma_2' B_2(x)$. Recalling the decomposition of $g(x)$ and that $\mathbf{B}(x)$ is an orthonormal basis with the normalization factor of

L^{-1} and $\|B_2\|_\infty = O(L^{-1/2})$, we get

$$L^{-1}|\gamma_2^* - \gamma_2'| = \left| \int_0^1 (g^*(x) - g(x))B_2(x)dx \right| \leq C_4 C_g L^{-5/2}.$$

Thus we have $|\gamma_2^* - \gamma_2'| \leq C_4 C_g L^{-3/2}$ and

$$\|(\gamma_2^* - \gamma_2')B_2\|_\infty \leq C_5 C_g L^{-2}. \quad (\text{S.50})$$

Note that C_3 , C_4 , and C_5 are independent of the specific function. Hence the desired result follows from (S.49) and (S.50).

Finally we consider

$$\Omega_1 = \int_0^1 \mathbf{B}'_0(z)(\mathbf{B}'_0(z))^T dz, \quad \Omega_2 = \int_0^1 \mathbf{B}''_0(z)(\mathbf{B}''_0(z))^T dz, \quad \text{and } \mathbf{B}_0(Z_1) - \mathbb{E}\{\mathbf{B}_0(Z_1)\}.$$

We demonstrate that both Ω_1 and Ω_2 does not necessarily have desirable properties for theoretical analysis. This conclusion also applies to $\mathbf{B}_0(Z_1) - \mathbb{E}\{\mathbf{B}_0(Z_1)\}$.

Take a three times continuously differentiable function $g(z)$. Then Corollary 6.26 in Schumaker (2007) implies that for some $\gamma \in R^L$,

$$\begin{aligned} \|g - \gamma^T \mathbf{B}_0\| &\leq C_1 L^{-3} \sum_{j=0}^3 \|g^{(j)}\|, \\ \|g' - \gamma^T \mathbf{B}'_0\| &\leq C_2 L^{-2} \sum_{j=0}^3 \|g^{(j)}\|, \\ \|g'' - \gamma^T \mathbf{B}''_0\| &\leq C_3 L^{-1} \sum_{j=0}^3 \|g^{(j)}\|. \end{aligned}$$

where C_1 , C_2 , and C_3 are independent of $g(z)$.

Taking $g(z) = \sin(2\pi Rz)$ with $R \rightarrow \infty$ and $R^3/L \rightarrow 0$, we have from the above three inequalities that

$$\begin{aligned} \|g\| &\sim 1, \quad \|g'\| \sim R, \quad \|g''\| \sim R^2, \\ \gamma^T \Omega_0 \gamma &\sim 1, \quad (\gamma^T \Omega_1 \gamma)^{1/2} \sim R, \quad (\gamma^T \Omega_2 \gamma)^{1/2} \sim R^2. \end{aligned}$$

These and (S.47) imply that Ω_1 and Ω_2 have eigenvalues $\tilde{\lambda}_1$ and $\tilde{\lambda}_2$ satisfying $\tilde{\lambda}_j L \rightarrow \infty$ ($j = 1, 2$), respectively. This contradicts some critical assumptions in some papers.

To consider $\mathbf{B}_0(Z_1) - \mathbb{E}\{\mathbf{B}_0(Z_1)\}$, we note the following equations.

$$\sum_{j=1}^L \tau_j = 1 \quad \text{and} \quad \begin{pmatrix} B_{02}(Z_1) - \mathbb{E}\{B_{02}(Z_1)\} \\ \vdots \\ B_{0L}(Z_1) - \mathbb{E}\{B_{0L}(Z_1)\} \end{pmatrix} = D\mathbf{B}_0(Z_1), \quad (\text{S.51})$$

where $\tau_j = E\{B_{0j}(Z_1)\}$ and the $(L-1) \times L$ matrix D is defined by

$$D = (0 I_{L-1}) - \begin{pmatrix} \tau_2 & \cdots & \tau_2 \\ \dots & \dots & \dots \\ \tau_L & \cdots & \tau_L \end{pmatrix}.$$

When Z_1 has a bounded density function, $\tau_j \sim 1/L$ uniformly in j and we have

$$\mathbf{i}_{L-1}^T D = (\tau_1 - 1, \tau_1, \dots, \tau_1) \text{ and } |D^T \mathbf{i}_{L-1}| \sim 1$$

for $\mathbf{i}_{L-1} = (1, \dots, 1)^T \in R^{L-1}$. This means

$$\lambda_{\min}(DD^T) = O(L^{-1}) \text{ and } \lambda_{\min}(D\Omega_0 D^T) = O(L^{-2}).$$

This implies that the basis in (S.51) is not suitable for additive models for this poor eigenvalue property. That is why we have introduced another basis.

S.4 A real application

In this section, the usefulness of our methods is illustrated via a real dataset, available from <http://www.res.org.uk>. This dataset includes a variety of variables describing more than 100 Japanese industrial chemical firms listed on the Tokyo stock exchange. The goal is to investigate the relationship between shareholder concentration and several indices for managerial moral hazard in the form of expenditure with scope for private benefit, which has been studied previously by Yafeh and Yosha (2003), Horowitz and Lee (2005) and Lian (2012). The response variable Y is the general sales and administrative expenses deflated by sales, which is one of five measures of activities with a scope for managerial moral hazard given by Yafeh and Yosha (2003). We consider the additive regression model with covariates, X_1 : log(assets), X_2 : the age of the firm, X_3 : leverage (ratio of debt to total assets), X_4 : profit (variance of operating profitability of firms between 1977 and 1986), X_5 : TOPTEN (the percentage of ownership held by the 10 largest shareholders), and X_6 : share (share of the largest creditor in total debt). All covariates are normalized into the range of $[0,1]$ via a linear transformation. Note that only 114 firms are included in our analysis because of the missing covariates. The model selection results based on AWG-Lasso+HDIC and AWG-Lasso+HDIC_{II} with $\tau = 0.5$ are summarized in Table S.1, in which L_i , NL_i , and $NS_i, i = 1, \dots, 6$, are defined as in Example 2 of Section 4 with the number of replications set to 1. It is shown in Table

Table S.1: $(L_i, NL_i, NS_i), i = 1, \dots, 6$ in the data on Japanese industrial chemical firms

	$(n, p) = (114, 6)$ with $\tau = 0.5$					
	(L_1, NL_1, NS_1)	(L_2, NL_2, NS_2)	(L_3, NL_3, NS_3)	(L_4, NL_4, NS_4)	(L_5, NL_5, NS_5)	(L_6, NL_6, NS_6)
AWG-Lasso+HDIC	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)
AWG-Lasso+HDIC _{II}	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(1.0, 0.0, 0.0)	(0.0, 1.0, 0.0)	(1.0, 0.0, 0.0)	(1.0, 0.0, 0.0)

S.1 that all covariates are survived after model selection, and the firm size $[\log(\text{assets})]$ and firm age [the age of the firm] are identified as nonlinear covariates, in accordance with the findings given by Horowitz and Lee (2005).

We now evaluate the performance of these two methods in high-dimensional situations based on the result obtained from the above analysis. To this aim, we artificially generate irrelevant covariates $X_j, j = 7, \dots, 100$, where X_7, \dots, X_{10} are i.i.d. $N(0, (0.15)^2)$ distributed, X_{11}, \dots, X_{100} are i.i.d. $U(0, 0.5)$ distributed, and (X_7, \dots, X_{10}) and (X_{11}, \dots, X_{100}) are independent. In view of (S.1), we also transform (X_7, \dots, X_{10}) into the range $[0, 1]$ using $(X_i - \min_{i \in \{7, \dots, 10\}} X_i) / (\max_{i \in \{7, \dots, 10\}} X_i - \min_{i \in \{7, \dots, 10\}} X_i)$, $i = 7, \dots, 10$. Since the above analysis suggests that all X_1, \dots, X_6 are relevant, the TNR is defined by

$$\frac{\sum_{j=7}^{100} I_{\{X_j \text{ is not selected}\}}}{94}.$$

In order to alleviate the overfitting problem of AWG-Lasso+HDIC and AWG-Lasso+HDIC_{II}, we also increase the penalty q_n from 1 to 3. For $\tau = 0.5$ and $q_n = 1, 2$ and 3, the performance of the two methods in terms of L_i, NL_i , and $NS_i, i = 1, \dots, 6$, and TNR is documented in Table S.2. Table S.2 shows that for $q_n = 1$, X_1, \dots, X_6 are still identified as relevant variables although the structure identification results may slightly differ from the low-dimensional case considered previously. However, the TNR values of AWG-Lasso+HDIC and AWG-Lasso+HDIC_{II} are 0.72 and 0.35, suggesting that the methods suffer from an overfitting problem. When q_n increases to 2, the TNR values of the two methods increase substantially, in particular, for AWG-Lasso+HDIC_{II}. On the other hand, X_2 and X_5 are now identified as irrelevant variables, suggesting the potential problem of false negatives. The result for $q_n = 3$ is similar to that for $q_n = 2$.

In conclusion, we note that our methods perform reasonably well in discovering relevant variables and excluding irrelevant ones. The model/structure identification results, however, may sometimes be sensitive to the choice of q_n , which deserves a separate investigation.

Table S.2: $(L_i, NL_i, NS_i), i = 1, \dots, 6$ and TNR in the data on Japanese industrial chemical firms with artificial covariates

$(n, p, q_n) = (114, 100, 1)$ with $\tau = 0.5$						
	(L_1, NL_1, NS_1)	(L_2, NL_2, NS_2)	(L_3, NL_3, NS_3)	(L_4, NL_4, NS_4)	(L_5, NL_5, NS_5)	(L_6, NL_6, NS_6)
AWG-Lasso+HDIC	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(1.0, 0.0, 0.0)	(0.0, 1.0, 0.0)
AWG-Lasso+HDIC _{II}	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)
	(TNR)					
AWG-Lasso+HDIC	(0.72)					
AWG-Lasso+HDIC _{II}	(0.35)					
$(n, p, q_n) = (114, 100, 2)$ with $\tau = 0.5$						
	(L_1, NL_1, NS_1)	(L_2, NL_2, NS_2)	(L_3, NL_3, NS_3)	(L_4, NL_4, NS_4)	(L_5, NL_5, NS_5)	(L_6, NL_6, NS_6)
AWG-Lasso+HDIC	(0.0, 1.0, 0.0)	(0.0, 0.0, 1.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 0.0, 1.0)	(0.0, 1.0, 0.0)
AWG-Lasso+HDIC _{II}	(0.0, 1.0, 0.0)	(0.0, 0.0, 1.0)	(0.0, 1.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 0.0, 1.0)	(0.0, 1.0, 0.0)
	(TNR)					
AWG-Lasso+HDIC	(0.82)					
AWG-Lasso+HDIC _{II}	(0.82)					
$(n, p, q_n) = (114, 100, 3)$ with $\tau = 0.5$						
	(L_1, NL_1, NS_1)	(L_2, NL_2, NS_2)	(L_3, NL_3, NS_3)	(L_4, NL_4, NS_4)	(L_5, NL_5, NS_5)	(L_6, NL_6, NS_6)
AWG-Lasso+HDIC	(0.0, 1.0, 0.0)	(0.0, 0.0, 1.0)	(1.0, 0.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 0.0, 1.0)	(0.0, 1.0, 0.0)
AWG-Lasso+HDIC _{II}	(0.0, 1.0, 0.0)	(0.0, 0.0, 1.0)	(1.0, 0.0, 0.0)	(0.0, 1.0, 0.0)	(0.0, 0.0, 1.0)	(0.0, 1.0, 0.0)
	(TNR)					
AWG-Lasso+HDIC	(0.82)					
AWG-Lasso+HDIC _{II}	(0.82)					

References

- [1] P. BÜHLMANN AND S. VAN DE GEER. (2011). *Statistics for High-Dimensional Data: Methods Theory and Applications*. Springer, New York, Dordrecht, Heidelberg, London.
- [2] J. FAN, Y. FAN, AND E. BARUT. (2014). Adaptive robust variable selection. *Ann. Statist.*, **42**, 324–351.
- [3] HOROWITZ, J. L., AND LEE, S. (2005). Nonparametric estimation of an additive quantile regression model. *J. Amer. Statist. Assoc.*, **100**, 1238–1249.
- [4] H. LIAN. (2012). Semiparametric estimation of additive quantile regression models by twofold penalty. *Journal of Business & Economic Statistics*, **30**, 337–350.
- [5] L. L. SCHUMAKER. (2007). *Spline Functions: Basic Theory* 3rd ed. Cambridge University Press, Cambridge, 2007.

- [6] B. SHERWOOD AND L. WANG. (2016). Partially linear additive quantile regression in ultra-high dimension. *Ann. Statist.*, **44**, 288–317.
- [7] Y. TANG, X. SONG, H. J. WANG, AND Z. ZHU. (2013). Variable selection in high-dimensional quantile varying coefficient models. *J. Multivariate Anal.*, **122**, 115–132.
- [8] S. VAN DE GEER. (2000). Empirical Processes in M-estimation. Cambridge University Press, Cambridge.
- [9] YAFEH, Y., AND YOSHA, O. (2003). Large Shareholders and Banks: Who Monitors and How? *The Economic Journal*, **113**, 128–146.