

Sparse Estimation in Semi-parametric Finite Mixture of Varying Coefficient Regression Models

Abbas Khalili^{1,*}, Farhad Shokoochi^{2,**}, Masoud Asgharian^{3,***}, and Shili Lin^{4,****}

^{1,3}Department of Mathematics and Statistics, McGill University, Montreal, QC, Canada

²Department of Mathematical Sciences, University of Nevada Las Vegas, Las Vegas, NV 89154, USA

⁴Department of Statistics, Ohio State University, Columbus, OH 43210, USA

**email*: abbas.khalili@mcgill.ca

***email*: farhad.shokoochi@unlv.edu

****email*: masoud.asgharian2@mcgill.ca

*****email*: shili@stat.osu.edu

SUMMARY: Finite mixture of regressions (FMR) are commonly used to model heterogeneous effects of covariates on a response variable in settings where there are unknown underlying subpopulations. FMRS, however, cannot accommodate situations where covariates' effects also vary according to an "index" variable—known as finite mixture of varying coefficient regression (FM-VCR). Although complex, this situation occurs in real data applications: the osteocalcin (OCN) data analyzed in this manuscript presents a heterogeneous relationship where the effect of a genetic variant on OCN in each hidden subpopulation varies over time. Oftentimes, the number of covariates with varying coefficients also presents a challenge: in the OCN study, genetic variants on the same chromosome are considered jointly. The relative proportions of hidden subpopulations may also change over time. Nevertheless, existing methods cannot provide suitable solutions for accommodating all these features in real data applications. To fill this gap, we develop statistical methodologies based on regularized local-kernel likelihood for simultaneous parameter estimation and variable selection in sparse FM-VCR models. We study large-sample properties of the proposed methods. We then carry out a simulation study to evaluate the performance of various penalties adopted for our regularized approach and ascertain the ability of a BIC-type criterion for estimating the number of subpopulations. Finally, we applied the FM-VCR model to analyze the OCN data and identified several covariates, including genetic variants, that have age-dependent effects on osteocalcin.

KEYWORDS: finite mixture of regressions, local-kernel likelihood, non-parametric models, penalized likelihood

1. Introduction

Finite mixture of regression (FMR) models ([McLachlan and Peel, 2000](#)) are commonly used to accommodate heterogeneous effects of covariates $\mathbf{X} = (X_1, \dots, X_d)$ on a response variable Y when the population under study is believed to consist of multiple hidden subpopulations. While FMRS can successfully capture such heterogeneity, they fall short if the effects of X_j 's on Y also vary as functions of an index variable, U , such as time or location. The mixing proportions, representing the sizes of hidden subpopulations, may also change as functions of U . It becomes even more challenging when only a few covariates among a large set have significant effects on Y . The data from the osteocalcin (OCN) study under our consideration presents such challenges on all fronts: there are many genetic and non-genetic covariates (X_j 's); there may be heterogeneous relationships, also varying over age (U), between the response variable (OCN, Y) and only a subset of the covariates; and the relative proportions of hidden subpopulations may also change over age.

An extension of FMRS, semi-parametric finite mixture of varying coefficient regressions (FM-VCR), was further introduced to account for heterogeneous varying covariates' effects ([Xiang et al., 2019](#)). These models facilitate the use of varying coefficient regressions, as functions of U , in studying the relationship between Y and X_j 's in a heterogeneous population. The functional forms of varying coefficients are seldom known and left unspecified. [McLachlan and Peel \(2000\)](#) and [Xiang et al. \(2019\)](#), discuss statistical inference and applications of various special cases of FM-VCRs in applied sciences and machine learning.

In the OCN study, the investigators collected data on environmental and genetic factors to study whether and how osteocalcin (the phenotype) is affected by a subset of these factors ([Liao et al., 2014](#)). Several genetic variants in Chromosome 7 (Chr7) were implicated ([Zhang, 2017](#)), and the gene harboring these variants was in fact linked to bone morphogenetic protein (BMP) ([Harada et al., 2003](#)). Population heterogeneity has not, however, been considered in

these studies. Inspired by recent studies demonstrating that relationships between OCN and some underlying factors are not the same across the population (Liu et al., 2015), we are interested in exploring whether there exist subpopulations, where the relationships between OCN and the factors of interest, especially the single-nucleotide polymorphisms (SNPs) in Chr7, are different across the subpopulations. Further, genetic effects on a phenotype may change over time (Zhang, 2017), which may well be the case for genes related to BMP. Finally, it is known that aging leads to bone density loss both in men and women (Demontiero et al., 2012); the mixing proportions of the components may well vary with age. Thus, a reanalysis of the data to accommodate population heterogeneity, varying subpopulation relative proportions, and varying genetic effects over time is warranted.

In this paper, motivated by the challenges in analyzing the OCN data, we study methodologies for sparse estimation in FM-VCR models. To the best of our knowledge, the existing literature (Xiang et al., 2019) mainly focuses on maximum likelihood estimation (MLE), while variable selection problems are largely understudied. In FM-VCRs, even with $d = 43$ in the OCN study, the dimension of the parameter space is large enough to render classical methods for variable selection almost impractical. Thus, we develop new results based on regularized local-kernel likelihood methods, demonstrating that regularized estimation in FM-VCRs consistently estimates the model parameters while recovering its sparse structure.

Since the seminal work of Hastie and Tibshirani (1993), varying coefficient regressions (VCR) and their extensions to mixture models (Xiang et al., 2019) have attracted much attention in the statistics literature. The parameters in these models can change as smooth non-parametric functions of an index variable such as time or location, which results in reducing modeling bias while avoiding the curse of dimensionality (Fan and Zhang, 1999). Due to the non-parametric nature of the models, parameter estimation requires careful consideration. Several estimation techniques are available in the literature, ranging from

the local kernel (Huang et al., 2018), to splines and basis functions (Hastie and Tibshirani, 1993) and local polynomial approximations (Fan and Gijbels, 1996). Another line of research particularly in VCRs has focused on variable selection, thanks to the LASSO (Tibshirani, 1996), SCAD (Fan and Li, 2001), adaptive LASSO (AdpLASSO) (Zou, 2006), MCP (Zhang, 2010), and group LASSO (Yuan and Lin, 2006). Among others, Wang and Xia (2009) studied variable selection using group LASSO/AdpLASSO and basis expansion or local-kernel methods. Wei et al. (2011) and others studied the problem in high dimensions using direct penalization or two-stage screening and penalization.

Despite the surge of research on variable selection in VCRs, the problem in FM-VCRs under our consideration has not been studied. The works reviewed by Xiang et al. (2019) focus on MLE in FM-VCRs and their special cases. One could use best-subset selection methods such as the AIC, BIC, or their variations to perform variable selection based on the MLE. Such techniques, however, require intensive computations as there are potentially $2^{C \times d} > 10^6$ submodels to be examined for selecting a sparse 2-components FM-VCR even with only $d = 10$ covariates. Variable selection based on regularization techniques such as LASSO is mainly studied in FMRS (Khalili and Chen, 2007; Städler et al., 2010; Shokoochi et al., 2019), without allowing for varying coefficients. The challenges of sparse estimation in FM-VCRs are the non-parametric nature of varying regression coefficients, mixing probabilities, and dispersion parameters, and often many covariates where only a handful are significant in the model.

In this paper, we develop computationally efficient penalized local-kernel likelihood methods for sparse estimation in FM-VCRs. We establish consistency in estimation and variable selection and oracle properties of the proposed estimators. We develop a modified EM algorithm (Dempster et al., 1977) for the numerical implementation of the methods. We evaluate the finite-sample performance of the methods via simulation and analyze the OCN data.

2. Sparse FM-VCR models

Let Y be a real-valued response variable, \mathbf{X} be a d -dimensional vector of covariates, and U be an index variable. In a population with $C > 1$ hidden subpopulations, we are interested in the conditional distribution of $Y|(\mathbf{X} = \mathbf{x}, U = u)$, for $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$ and $u \in \mathcal{U} \subset \mathbb{R}$.

Let \mathcal{C} be a discrete latent variable taking values $1, \dots, C$. For any $\mathbf{x} \in \mathcal{X}, u \in \mathcal{U}$, we denote

$$0 < \pi_j(u) = \Pr(\mathcal{C} = j | \mathbf{X} = \mathbf{x}, U = u), \quad j = 1, \dots, C, \quad (1)$$

as non-parametric functions of u , and $\sum_{j=1}^C \pi_j(u) = 1$. This setting differs from the mixture-of-experts (ME) and their variations (McLachlan and Peel, 2000) including finite mixtures with concomitant variables (Dayton and Macready, 1988; Frühwirth-Schnatter et al., 2018), wherein each π_j is modeled as a parametric function of u and possibly other covariates \mathbf{x} .

Suppose the conditional density (mass) function of $Y|(\mathbf{X} = \mathbf{x}, U = u, \mathcal{C} = j)$ is given by

$$f(y; \theta_j(\mathbf{x}, u), \phi_j(u)), \quad \text{for any } y \in \mathcal{Y} \subset \mathbb{R}, \quad (2)$$

where f belongs to a known parametric family $\mathcal{G} = \{f(y; \theta, \phi) : (\theta, \phi) \in \Theta \times (0, \infty), \Theta \subset \mathbb{R}\}$ with respect to a σ -finite measure ν , and $\phi > 0$ is a dispersion parameter. We assume $\theta_j(\mathbf{x}, u) = g(\mathbf{x}^\top \boldsymbol{\beta}_j(u))$, for a known link function g and the regression coefficients $\boldsymbol{\beta}_j^\top(u) = (\beta_{j1}(u), \dots, \beta_{jd}(u))$, where $(\phi_j(u), \beta_{jl}(u))$ are also unknown non-parametric functions of u . A well-known example of f in (2) is the Gaussian density with the mean $\theta_j(\mathbf{x}, u) = E\{Y|(\mathbf{x}, u, j)\} = \mathbf{x}^\top \boldsymbol{\beta}_j(u)$ which corresponds to the identity link function $g(\eta) = \eta$, and variance $\text{Var}\{Y|(\mathbf{x}, u, j)\} = \phi_j(u)$; see Web Appendix F for more details. Also note that each $\boldsymbol{\beta}_j(u)$ and $\phi_j(u)$ is a function of u rather than a constant as in the ME models.

By putting (1) and (2) together, the conditional density (mass) function of $Y|(\mathbf{X} = \mathbf{x}, U = u)$ in an FM-VCR model with order C is given by

$$f_C^*(y | \boldsymbol{\psi}(u), \mathbf{x}) = \sum_{j=1}^C \pi_j(u) f(y; \theta_j(\mathbf{x}, u), \phi_j(u)), \quad y \in \mathcal{Y} \quad (3)$$

where $\boldsymbol{\psi}(u) = (\boldsymbol{\pi}^\top(u), \boldsymbol{\phi}^\top(u), \boldsymbol{\beta}_1^\top(u), \dots, \boldsymbol{\beta}_C^\top(u))^\top \in \mathbb{R}^p$, with $p = Cd + 2C - 1$, and the sub-vectors $\boldsymbol{\phi}(u) = (\phi_1(u), \dots, \phi_C(u))^\top \in \mathbb{R}^C$, $\boldsymbol{\pi}(u) = (\pi_1(u), \dots, \pi_{C-1}(u))^\top \in \mathbb{R}^{(C-1)}$.

Identifiability is essential in mixture models; if for any $C_1 > 1, C_2 > 1$ and given (u, \mathbf{x}) , we have $f_{C_1}^*(y|\boldsymbol{\psi}_1(u), \mathbf{x}) = f_{C_2}^*(y|\boldsymbol{\psi}_2(u), \mathbf{x})$, for all $y \in \mathcal{Y}$, then we must have $C_1 = C_2$ and $\boldsymbol{\psi}_1(u) = \boldsymbol{\psi}_2(u)$ (up to a mixture component permutation). Sufficient conditions for identifiability of (3) are discussed in Theorem 1 of Huang et al. (2018), and the conditions are given in Web Appendix A. We assume that the models considered here are identifiable.

In an FM-VCR, even for moderate values of C and d , the number $p = Cd + (2C - 1)$ of non-parametric functions to be estimated becomes very large compared to a typical sample size n . In our real data application for Chr7, after an initial covariate screening, we have 36 SNPs plus 7 extra covariates ($d = 43$), with $1 \leq C \leq 5$ resulting in potentially $44 \leq p \leq 224$ non-parametric functions to be estimated based on a sample size of $n = 1704$. However, not all the covariates have significant effects on the response variable of interest, suggesting a sparse model for the data. Thus, we assume that the true FM-VCR model is *sparse*, that is, for every $j = 1, \dots, C$, there exists an integer $1 \leq d_j^0 < d$ such that

$$0 < E[\{\beta_{jl}^0(U)\}^2] < \infty, \quad \text{for all } l = 1, \dots, d_j^0, \quad \text{and} \quad (4)$$

$$E[\{\beta_{jl}^0(U)\}^2] = 0, \quad \text{for all } l = d_j^0 + 1, \dots, d, \quad (5)$$

where E is with respect to U , and d_j^0 is the number of non-zero $\beta_{jl}^0(\cdot)$'s in component j . By (4)-(5), the total number of non-parametric functions is $p_0 = \sum_{j=1}^C d_j^0 + (2C - 1)$ and presumably much smaller than p . More discussion on Condition (5) is given in Remark 1 of Web Appendix E. Next, we discuss estimation and variable selection in sparse FM-VCRs.

3. Simultaneous estimation and variable selection

Let (u_i, \mathbf{x}_i, y_i) , $i = 1, \dots, n$, be the observations based on a random sample from the FM-VCR model (3). The (conditional) log-likelihood function is given by

$$\mathcal{L} = \sum_{i=1}^n \log\{f_C^*(y_i|\boldsymbol{\psi}(u_i), \mathbf{x}_i)\} = \sum_{i=1}^n \log\left\{\sum_{j=1}^C \pi_j(u_i) f(y_i; \boldsymbol{\theta}_j(\mathbf{x}_i, u_i), \phi_j(u_i))\right\}, \quad (6)$$

where $\boldsymbol{\theta}_j(\mathbf{x}_i, u_i) = g(\mathbf{x}_i^\top \boldsymbol{\beta}_j(u_i))$. Since $(\pi_j(u), \phi_j(u), \boldsymbol{\beta}_j(u))$ are unknown non-parametric functions, the parameter space has infinite dimension and \mathcal{L} is thus intractable. Several

techniques are available in the literature ([Fan and Gigbels, 1996](#)) for non-parametric function estimation, of which the most popular are local-kernel methods, splines and local polynomial approximations. [Silverman \(1984\)](#) shows that spline smoothing corresponds approximately to smoothing by a kernel method with bandwidth depending on the local density of design points. In the following, we use the local-kernel method where each functional parameter in the FM-VCR model is locally approximated by constants. As a result, the locally approximated model becomes a standard finite mixture of regression which in turn makes implementation of our penalization method via the EM algorithm easier ([Huang et al., 2013](#)).

Local log-likelihood. For any $u \in \mathcal{U}$, we consider the local-kernel log-likelihood

$$\ell_n(\boldsymbol{\psi}(u); h) = \sum_{i=1}^n \log \left\{ \sum_{j=1}^C \pi_j f(y_i; \theta_j(\mathbf{x}_i), \phi_j) \right\} K_h(u_i - u), \quad (7)$$

where $\theta_j(\mathbf{x}_i) = g(\mathbf{x}_i^\top \boldsymbol{\beta}_j)$, and $K_h(t) = h^{-1}K(t/h)$ is a kernel function with a band-width h . The locally constant vector of parameters is $\boldsymbol{\psi}(u) = (\boldsymbol{\pi}^\top, \boldsymbol{\phi}^\top, \boldsymbol{\beta}_1^\top, \dots, \boldsymbol{\beta}_C^\top) \in \mathbb{R}^p$, where $\boldsymbol{\pi}^\top = (\pi_1, \dots, \pi_{C-1})$, $\boldsymbol{\phi}^\top = (\phi_1, \dots, \phi_C)$, and $\boldsymbol{\beta}_j^\top = (\beta_{j1}, \dots, \beta_{jd})$, $j = 1, \dots, C$. The entries of $\boldsymbol{\psi}(u)$ are local-constant approximations of the functions $(\pi_j(u), \phi_j(u), \beta_{jl}(u))$, which depend on u , and for simplicity, we suppress u in the notation but keep it for $\boldsymbol{\psi}(u)$.

For any $u \in \mathcal{U}$, the maximum local-kernel log-likelihood estimator (MLLE) of $\boldsymbol{\psi}(u)$ is defined as the maximizer of $\ell_n(\boldsymbol{\psi}(u); h)$, and is denoted by $\check{\boldsymbol{\psi}}_n(u)$. In practice, the estimation is usually done at the observed points u_1, \dots, u_n ; more discussion is given at the end of this section. In Proposition 1 of [Web Appendix D](#), we establish estimation consistency of $\check{\boldsymbol{\psi}}_n(u)$. It is, however, well-known that similar to the MLE in parametric regression, $\check{\boldsymbol{\psi}}_n(u)$ does not provide a sparse FM-VCR model as postulated in (5), which is the main focus of our research. Therefore, we propose a regularization method that yields fitted sparse FM-VCR models.

Penalized local log-likelihood. For $j = 1, \dots, C$ and $t = 1, \dots, n$, let $\boldsymbol{\beta}_{j,t} = (\beta_{j1,t}, \dots, \beta_{jd,t})^\top$ be the local constant approximation of the vector $\boldsymbol{\beta}_j(u)$ at u_t 's. We denote the $(n \times d)$ matrices $\mathbf{B}_j = (\boldsymbol{\beta}_{j,1}, \dots, \boldsymbol{\beta}_{j,n})^\top$. Similarly, let $\boldsymbol{\pi}_j^\top = (\pi_{j,1}, \dots, \pi_{j,n})$ and $\boldsymbol{\phi}_j^\top = (\phi_{j,1}, \dots, \phi_{j,n})$ be the

local constant approximations of $\pi_j(u)$ and $\phi_j(u)$ at u_t 's, respectively. Also, denote the $(n \times p)$ matrix $\Psi = (\boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_{C-1}, \boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_C, \mathbf{B}_1, \dots, \mathbf{B}_C)$, where $p = Cd + 2C - 1$.

Using (7), the corresponding (total) local-kernel log-likelihood is given by

$$L_n(\Psi; h) = \sum_{t=1}^n \ell_n(\boldsymbol{\psi}(u_t); h) = \sum_{t,i=1}^n \log \left\{ \sum_{j=1}^C \pi_{j,t} f(y_i; \theta_{j,t}(\mathbf{x}_i), \phi_{j,t}) \right\} K_h(u_i - u_t), \quad (8)$$

and $\theta_{j,t}(\mathbf{x}_i) = g(\mathbf{x}_i^\top \boldsymbol{\beta}_{j,t})$. We estimate Ψ by maximizing the penalized local log-likelihood

$$\tilde{L}_n(\Psi; \boldsymbol{\lambda}, h) = L_n(\Psi; h) - \mathbb{P}_n(\Psi; \boldsymbol{\lambda}), \quad (9)$$

$$\mathbb{P}_n(\Psi; \boldsymbol{\lambda}) = \sum_{j=1}^C \sum_{l=1}^d p_n(\|\mathbf{b}_{jl}\|/\sqrt{n}; \lambda_j), \quad (10)$$

where p_n is a penalty function with tuning parameters $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_C)^\top$ that control the level of penalization on $\|\mathbf{b}_{jl}\|/\sqrt{n}$. Here, $\mathbf{b}_{jl} = (\beta_{jl,1}, \dots, \beta_{jl,n})^\top = (\beta_{jl}(u_1), \dots, \beta_{jl}(u_n))^\top$ is the l^{th} column of the matrix \mathbf{B}_j , $j = 1, \dots, C$, and $\|\mathbf{b}_{jl}\|^2 = \sum_{t=1}^n \beta_{jl}^2(u_t)$. Examples of p_n are the LASSO, AdpLASSO, SCAD, and MCP (Web Appendix C). The group penalization in (10) (Yuan and Lin, 2006) enforces zero estimates of some of $\|\mathbf{b}_{jl}\|$ and thus the vectors \mathbf{b}_{jl} . By Condition (RC.2) on the density $m(u)$ of U , u_1, \dots, u_n are dense in the support \mathcal{U} (Janson, 1987), and thus it suffices to perform estimation at only u_i 's rather than on the entire \mathcal{U} .

Given $(\boldsymbol{\lambda}, h)$, the maximum penalized local log-likelihood estimate (MPLLE) of Ψ is

$$\hat{\Psi}_n(\boldsymbol{\lambda}, h) \equiv \hat{\Psi}_n = \underset{\Psi}{\operatorname{argmax}} \tilde{L}_n(\Psi; \boldsymbol{\lambda}, h). \quad (11)$$

We induce zero estimates for some of the regression coefficients $\beta_{jl}(\cdot)$ by appropriate choices of the penalty p_n , and the tuning parameters $(\boldsymbol{\lambda}, h)$. Thus, the MPLLE performs simultaneous estimation and variable selection, resulting in a fitted sparse FM-VCR model.

4. Large-sample theory

To distinguish from earlier notation, while drawing the connection, for any $u \in \mathcal{U}$, let $\boldsymbol{\psi}^0(u)$ be the p -dimensional parameter vector in (3) corresponding to the true sparse FM-VCR model satisfying (4)-(5). In particular, for the observed points u_1, \dots, u_n , the $(n \times p)$ -dimensional parameter matrix in (8) corresponding to the true model is denoted by Ψ_n^0 with

its t^{th} row as $[\boldsymbol{\psi}^0(u_t)]^\top, t = 1, \dots, n$. Without loss of generality, we assume the partitioning $\boldsymbol{\psi}^0(u) = (\boldsymbol{\psi}_1^0(u), \boldsymbol{\psi}_2^0(u))$ such that $\boldsymbol{\psi}_2^0(u)$ contains only those $\beta_{jl}^0(\cdot)$'s that satisfy the sparsity Condition (5), and $\boldsymbol{\psi}_1^0(u)$ consists of the non-zero functions $\beta_{jl}^0(\cdot)$ among other parameters $(\pi_j^0(u), \phi_j^0(u))$. By (5), we have $E\{\|\boldsymbol{\psi}_2^0(U)\|^2\} = 0$. Thus, Conditions (C1) and (RC.2) on $\beta_{jl}^0(\cdot)$'s and the density $m(u)$ of U (Web Appendix A and Web Appendix B) imply that $\boldsymbol{\psi}_2^0(u) = \mathbf{0}$, uniformly in $u \in \mathcal{U}$. This results in $\boldsymbol{\Psi}_n^0 = (\boldsymbol{\Psi}_{n1}^0, \boldsymbol{\Psi}_{n2}^0)$ such that $\boldsymbol{\Psi}_{n2}^0$ contains all the zero regression functions, and $\boldsymbol{\Psi}_{n1}^0 = \{[\boldsymbol{\psi}_1^0(u_t)]^\top : t = 1, \dots, n\}$ contains the nonzero regression functions and other parameters. Denote $l(\boldsymbol{\psi}^0(u), \mathbf{x}, y) = \log\{f_C^*(y|\boldsymbol{\psi}^0(u), \mathbf{x})\}$, $y \in \mathcal{Y}$, $\mathbf{x} \in \mathcal{X}$. For $t = 1, \dots, n$, define $l'(\boldsymbol{\psi}^0(u_t), \mathbf{x}, y)$ and $l''(\boldsymbol{\psi}^0(u_t), \mathbf{x}, y)$, respectively, as the gradient and Hessian of $l(\boldsymbol{\psi}, \mathbf{x}, y)$ with respect to $\boldsymbol{\psi}$ and evaluated at $\boldsymbol{\psi} = \boldsymbol{\psi}^0(u_t)$. Also, let

$$\mathbf{I}(u_t) = E_{(Y, \mathbf{X})|U=u_t} [l'(\boldsymbol{\psi}^0(U), \mathbf{X}, Y)l'^\top(\boldsymbol{\psi}^0(U), \mathbf{X}, Y)|U = u_t)] \quad (12)$$

$$\Delta(u; u_t) = \int_{\mathcal{Y}} \int_{\mathcal{X}} l'(\boldsymbol{\psi}^0(u_t), \mathbf{x}, y) f_C^*(y|\boldsymbol{\psi}^0(u), \mathbf{x}) g(\mathbf{x}|u) d\mathbf{x} dy, \quad (13)$$

and $E_{(Y, \mathbf{X})|U=u_t}$ is the expectation with respect to the distribution of (\mathbf{X}, Y) given $U = u_t$.

With regard to the penalty function p_n , we denote the quantity

$$r_{1n} = \max \{h_n^{1/2}|p'_n(\theta_{jl}^0; \lambda_n)|/n^{3/2} : \theta_{jl}^0 = \sqrt{E_0\{\beta_{jl}^0(U)\}^2} \neq 0, 1 \leq j \leq C, 1 \leq l \leq d_j^0\} \quad (14)$$

where $p'_n(\theta_{jl}^0; \lambda_n)$ is the first derivative of $p_n(\theta_{jl}; \lambda_n)$ evaluated at $\theta_{jl} = \theta_{jl}^0 \neq 0$. Also, denote

$$\mathcal{V}_0 = \int_{\mathcal{U}} K^2(t) dt, \quad \mathcal{K}_2 = \int_{\mathcal{U}} t^2 K(t) dt. \quad (15)$$

Regularity Conditions (RC.1)-(RC.5) on the mixture density f_C^* , the conditional density $g(\mathbf{x}|u)$ of $\mathbf{X}|U = u$, the density $m(u)$ of U , and Conditions (KC.1)-(KC.2) and (PC.1)-(PC.3) on the kernel K and the penalty p_n , respectively, are all given in Web Appendix B. The proofs of the following lemma and theorems are given in Web Appendix D.

THEOREM 1: *(Point-wise estimation consistency) Suppose that Conditions (RC.1)-(RC.4) hold, the kernel K satisfies Conditions (KC.1)-(KC.2), and (λ_n, p_n, h_n) satisfy Conditions (PC.1)-(PC.2). Then, there exists a local maximizer $\widehat{\boldsymbol{\Psi}}_n$ of \widetilde{L}_n in (9) that $n^{-1}\|\widehat{\boldsymbol{\Psi}}_n - \boldsymbol{\Psi}_n^0\|_F^2 =$*

$n^{-1} \sum_{t=1}^n \|\widehat{\boldsymbol{\psi}}_n(U_t) - \boldsymbol{\psi}^0(U_t)\|^2 = O_p\{(1 + r_{1n})^2(nh_n)^{-1}\}$, where r_{1n} is given in (14), $\widehat{\boldsymbol{\psi}}_n(U_t)$ and $\boldsymbol{\psi}^0(U_t)$ are the t^{th} rows of the matrices $\widehat{\boldsymbol{\Psi}}_n$ and $\boldsymbol{\Psi}_n^0$, and $\|\cdot\|_F$ is the Frobenius norm.

By Theorem 1, if $r_{1n} = O(1)$, $\widehat{\boldsymbol{\Psi}}_n$ achieves the point-wise consistency rate $\{nh_n\}^{-1/2}$ in estimating $\boldsymbol{\Psi}_n^0$, a property shared by the MLLE (Proposition 1, Web Appendix D). This result clearly depends on the choice of (λ_n, p_n, h_n) . For the LASSO, $\lambda_n = O(\{nh_n\}^{-1/2})$ suffices. For SCAD and MCP, as long as $(\lambda_n, h_n) \rightarrow 0$ as $n \rightarrow \infty$, the desired rate is achieved since $r_{1n} = 0$. For AdpLASSO, with the (possibly random) weights w_{jl} 's, we require $\lambda_n \max_{1 \leq j \leq C, 1 \leq l \leq d_j^0} w_{jl} = o_p(\{nh_n\}^{-1/2})$. The practical choice of w_{jl} is given in Remark 3 of Web Appendix E. Theorem 1, however, does not imply the sparsity of $\widehat{\boldsymbol{\Psi}}_n$. We strengthen this result by establishing the sparsity and oracle property of the MPLLE, beginning with a lemma.

LEMMA 1: *Assume that the conditions of Theorem 1 are met, and that (λ_n, p_n, h_n) also satisfy Condition (PC.3). Then, for any $\sqrt{nh_n}$ -consistent MPLLE $\widehat{\boldsymbol{\Psi}}_n$ of $\boldsymbol{\Psi}_n^0$, as $n \rightarrow \infty$, $P\{\sum_{t=1}^n \widehat{\beta}_{jl}^2(U_t)/n\}^{1/2} = \|\widehat{\mathbf{b}}_{jl}/n\| = 0\} \rightarrow 1$, $j = 1, \dots, C$, $l = d_{j+1}^0, \dots, d$.*

By Lemma 1, with probability tending to one, the MPLLE of those $\beta_{jl}(u)$ that satisfy the sparsity property (5) are zero at the observed index variable: $\widehat{\beta}_{jl}(U_t) = 0, t = 1, \dots, n$ (point-wise sparsity). It is more appealing to establish a uniform sparsity. Recall the partitioning $\boldsymbol{\psi}^0(u) = (\boldsymbol{\psi}_1^0(u), \boldsymbol{\psi}_2^0(u))$, where $\boldsymbol{\psi}_2^0(u) = \mathbf{0}$ uniformly in $u \in \mathcal{U}$. We denote the oracle estimator $\widehat{\boldsymbol{\psi}}_{n1,orc}(u)$ as the MLLE of $\boldsymbol{\psi}_1^0(u)$, having known $\boldsymbol{\psi}_2^0(u) = \mathbf{0}$ a priori. For any $u \in \mathcal{U}$, let $\widehat{\boldsymbol{\psi}}_n(u) = (\widehat{\boldsymbol{\psi}}_{n1}(u), \widehat{\boldsymbol{\psi}}_{n2}(u))$ be the MPLLE of $\boldsymbol{\psi}^0(u)$, where its partitioning corresponds to that of $\boldsymbol{\psi}^0(u)$. Next, we describe the behavior of $\widehat{\boldsymbol{\psi}}_n(u)$ with respect to the oracle estimator.

THEOREM 2: *Under the conditions of Lemma 1, and Condition (RC.5), as $n \rightarrow \infty$,*

(i) *uniform sparsity:* $P\{\sup_{u \in \mathcal{U}} \|\widehat{\boldsymbol{\psi}}_{n2}(u)\| = 0\} \rightarrow 1$; (ii) *oracle property:* $\sup_{u \in \mathcal{U}} \|\widehat{\boldsymbol{\psi}}_{n1}(u) - \widehat{\boldsymbol{\psi}}_{n1,orc}(u)\| = o_p\{(\log h_n^{-1}/nh_n)^{1/2}\}$; (iii) *asymptotic normality:* for any $u \in \mathcal{U}$,

$$\sqrt{nh_n} \{[\mathbf{I}(u) + \mathbb{P}''_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)/n^2 m(u)](\widehat{\boldsymbol{\psi}}_{n1}(u) - \boldsymbol{\psi}_1^0(u)) - \mathcal{B}_n(u)\} \rightarrow N(\mathbf{0}, \mathcal{V}_0 \mathbf{I}(u)/m(u)),$$

where $\mathbf{I}(u)$ in (12) is positive definite, the bias $\mathcal{B}_n(u) = \frac{1}{2}\mathcal{K}_2\{\Delta''(u; u) + 2\Delta'(u; u)m'(u)/m(u)\}h_n^2 - \mathbb{P}'_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)/n^2m(u) + o_p(h_n^2)$, and $\boldsymbol{\theta}_1^0 = \{\theta_{jl}^0 = (E_0\{\beta_{jl}^0(U)\}^2)^{1/2} : 1 \leq j \leq C, 1 \leq l \leq d_j^0\}$.

Theorem 2(i) implies that the MPLLE of those $\beta_{jl}(\cdot)$'s that satisfy the sparsity property (5) are zero, that is $\widehat{\beta}_{jl}(u) = 0$ uniformly in $u \in \mathcal{U}$, for all $j = 1, \dots, C$ and $l = d_{j+1}^0, \dots, d$. Condition (PC.3) on (λ_n, p_n, h_n) guarantees sparsity of the MPLLE. For the LASSO, SCAD, and MCP, this condition requires $\sqrt{nh_n}\lambda_n \rightarrow \infty$, and for the AdpLASSO we need $\sqrt{nh_n}\lambda_n \min_{1 \leq j \leq K, d_j^0+1 \leq l \leq d} w_{jl} \rightarrow \infty$, as $n \rightarrow \infty$. Theorem 2(ii) and (iii) imply that the MPLLE asymptotically attains a similar information bound as the oracle estimator described above (Huang et al., 2018). Regarding (iii), for all the four penalties, we have $\mathbb{P}''_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n) = \mathbf{0}$. On the other hand, with respect to the bias $\mathcal{B}_n(u)$, for the SCAD and MCP, we have $\mathbb{P}'_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)/n^2 = \mathbf{0}$, and for the AdpLASSO, $\mathbb{P}'_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)/n^2 \sim \lambda_n \max_{1 \leq j \leq C, 1 \leq l \leq d_j^0} w_{jl}$, which tends to zero when scaled by $\sqrt{nh_n}$, as $n \rightarrow \infty$. For the LASSO, $\mathbb{P}'_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)/n^2 = \lambda_n$, which tends to ∞ once scaled by $\sqrt{nh_n}$, as required by Condition (PC.3) for sparsity. This behavior is also known in parametric regression, as it gives sparse estimators but does not achieve the oracle property. Finally, the asymptotic normality in (iii) is obtained from the oracle perspective that the true sparse structure of the model is known in advance. In practice, a sparse model is typically fitted based on the MPLLE. As such there are two sources of uncertainty that are caused by variable selection and parameter estimation. The asymptotic normality does not take the uncertainty due to variable selection into account, which is a topic of post-selection inference (PoSI, Berk et al. (2013)). More discussion is given in Remark 2 of Web Appendix E.

5. Computational strategies

The goal is to estimate $(\pi_j(u), \phi_j(u), \boldsymbol{\beta}_j(u))$, $j = 1, \dots, C$, at points u_1, \dots, u_n , by obtaining an approximate solution for (11). To avoid mixture label switching (Huang et al., 2013), we use a modified EM algorithm for estimation of each function simultaneously over u_1, \dots, u_n .

We view the observations $\{(u_i, \mathbf{x}_i, y_i) : i = 1, \dots, n\}$ as incomplete data, and introduce

the unobserved Bernoulli variables Z_{ij} to represent the membership of the i^{th} observation to the j^{th} mixture component, $j = 1, \dots, C$, such that $\Pr(Z_{ij} = 1 | u_i, \mathbf{x}_i, y_i) = \pi_j(u_i)$. The complete data are $\{(u_i, \mathbf{x}_i, y_i, \mathbf{z}_i) : i = 1, \dots, n\}$, where $\mathbf{z}_i = (z_{i1}, \dots, z_{iC})^\top$. We define the complete (total) local-kernel log-likelihood function as

$$L_n^c(\Psi; h) = \sum_{t=1}^n \sum_{i=1}^n \sum_{j=1}^C Z_{ij} \{ \log \pi_{j,t} + \log f(y_i; \theta_{j,t}(\mathbf{x}_i), \phi_{j,t}) \} K_h(u_i - u_t),$$

and the penalized complete local-kernel log-likelihood as $\tilde{L}_n^c(\Psi; h) = L_n^c(\Psi; h) - \mathbb{P}_n(\Psi; \lambda)$, where the penalty \mathbb{P}_n is given in (10). Due to the non-differentiability of $p_n(\theta; \lambda)$ at $\theta = 0$, as suggested by Fan and Li (2001), we use the local quadratic approximation (LQA)

$$\bar{p}_n(\theta; \theta^{(0)}, \lambda) = p_n(\theta^{(0)}; \lambda) + p'_n(\theta^{(0)}; \lambda)(\theta^2 - \theta^{2(0)})/(2\theta^{(0)}), \quad (16)$$

for p_n , where $\theta^{(0)}$ is an initial value. The LQA is used in the modified EM algorithm. Given $\Psi^{(m)}$, at the $(m+1)^{\text{th}}$ iteration, the algorithm proceeds as follows.

E-step. Since the Z_{ij} 's are unobservable, we compute the expectation of $\tilde{L}_n^c(\Psi; h)$ with respect to Z_{ij} conditional on $\{(u_i, \mathbf{x}_i, y_i) : i = 1, \dots, n\}$ and $\Psi^{(m)}$. For all $i = 1, \dots, n$, $j = 1, \dots, C$, it boils down to the computation of the conditional expectations

$$E(Z_{ij} | \Psi^{(m)}, u_i, \mathbf{x}_i, y_i) \equiv w_{ij}^{(m)} = \frac{\pi_j^{(m)}(u_i) f(y_i; \theta_j^{(m)}(\mathbf{x}_i, u_i), \phi_j^{(m)}(u_i))}{\sum_{j=1}^C \pi_j^{(m)}(u_i) f(y_i; \theta_j^{(m)}(\mathbf{x}_i, u_i), \phi_j^{(m)}(u_i))}.$$

M-step. Using the LQA (16) evaluated at $\theta = \|\mathbf{b}_{jl}\|_2/\sqrt{n}$ and $\theta^{(m)} = \|\mathbf{b}_{jl}^{(m)}\|_2/\sqrt{n}$, we maximize the objective function (up to some constants)

$$\tilde{Q}(\Psi; \Psi^{(m)}) = \sum_{j=1}^C \sum_{t=1}^n \left\{ \sum_{i=1}^n w_{ij}^{(m)} \{ \log \pi_{j,t} + \log f(y_i; \theta_{j,t}(\mathbf{x}_i), \phi_{j,t}) \} K_h(u_i - u_t) - \boldsymbol{\beta}_{j,t}^\top \boldsymbol{\Sigma}_j^{(m)} \boldsymbol{\beta}_{j,t} / 2 \right\}$$

with respect to Ψ , where $\theta_{j,t}(\mathbf{x}_i) = g(\mathbf{x}_i^\top \boldsymbol{\beta}_{j,t})$ with the vector $\boldsymbol{\beta}_{j,t} = (\beta_{j1}(u_t), \dots, \beta_{jd}(u_t))^\top$, the diagonal matrices $\boldsymbol{\Sigma}_j^{(m)} = \text{diag}\{\tau_{jl}^{(m)} : l = 1, \dots, d\}$ and $\tau_{jl}^{(m)} = p'_n(\theta^{(m)}; \lambda_j) / [n\theta^{(m)}]$.

The maximization of \tilde{Q} with respect to Ψ results in the following updates. The probabilities $\pi_j(u_t)$ and the vectors $\boldsymbol{\beta}_{j,t} = \boldsymbol{\beta}_j(u_t)$, $t = 1, \dots, n$, $j = 1, \dots, C$, are updated by

$$\pi_{j,t}^{(m+1)} = \pi_j^{(m+1)}(u_t) = \sum_{i=1}^n w_{ij}^{(m)} K_h(u_i - u_t) / \sum_{i=1}^n K_h(u_i - u_t),$$

$$\boldsymbol{\beta}_j^{(m+1)}(u_t) = \arg \max_{\boldsymbol{\beta}_{j,t} \in \mathbb{R}^d} \sum_{i=1}^n w_{ij}^{(m)} \{ \log f(y_i; \boldsymbol{\theta}_{j,t}(\mathbf{x}_i), \phi_{j,t}) \} K_h(u_i - u_t) - \boldsymbol{\beta}_{j,t}^\top \boldsymbol{\Sigma}_j^{(m)} \boldsymbol{\beta}_{j,t} / 2. \quad (17)$$

The dispersion parameters $\phi_{j,t} = \phi_j(u_t)$ are updated by solving the estimation equations

$$\sum_{i=1}^n w_{ij}^{(m)} \frac{\partial}{\partial \phi_{j,t}} \{ \log f(y_i; \boldsymbol{\theta}_{j,t}^{(m+1)}(\mathbf{x}_i), \phi_{j,t}) \} K_h(u_i - u_t) = 0. \quad (18)$$

To solve (17)-(18), depending on f , we may need to use the Newton-Raphson method.

Details of the algorithm for Gaussian and t -distribution f are given in [Web Appendix F](#).

Starting from an initial $\boldsymbol{\Psi}^{(0)}$ (discussed in [Web Appendix F.4](#)), the EM algorithm iterates until a convergence criterion is satisfied. We used the stopping rule $\|\boldsymbol{\Psi}^{(m+1)} - \boldsymbol{\Psi}^{(m)}\|_F < \epsilon$. We set the estimates of the regression functions $\beta_{jl}(u)$ to zero if $\|\mathbf{b}_{jl}^{(m+1)}/n\| < \delta$. In our simulation and the real data analysis, we chose $\epsilon = \delta = 10^{-4}$. We describe data-adaptive strategies for the selection of the band-width h and the tuning parameters λ_j 's, and also estimation of the mixture order C in [Web Appendix G](#).

[Table 1 about here.]

6. Simulation study

We assess the finite-sample performance of the methods via simulations. We considered Gaussian FM-VCRs with $C = 2$ and 3 components. The parameter settings for the model with $C = 2$ are given in [Table 1](#), and those for $C = 3$ are given in [Web Table 1](#).

Let CEZ = # Correctly Estimated Zero, CEN = # Correctly Estimated Non-zero, IEZ = # Incorrectly Estimated Zero, and IEN = # Incorrectly Estimated Non-zero $\beta_{jl}(\cdot)$'s. We define

$$\text{Sensitivity} = \text{CEN}/(\text{CEN} + \text{IEZ}) \quad \text{and} \quad \text{Specificity} = \text{CEZ}/(\text{CEZ} + \text{IEN}),$$

which assess the performance of the penalization method for variable selection. We measure the estimation errors of the proposed estimators using

$$L_2^2(\widehat{\boldsymbol{\beta}}_j) = n^{-1} \sum_{i=1}^n \sum_{l=0}^d (\widehat{\beta}_{jl}(u_i) - \beta_{jl}(u_i))^2, \quad L_2^2(\widehat{\boldsymbol{\sigma}}_j) = n^{-1} \sum_{i=1}^n (\widehat{\sigma}_j(u_i) - \sigma_j(u_i))^2, \\ L_2^2(\widehat{\boldsymbol{\pi}}_j) = n^{-1} \sum_{i=1}^n (\widehat{\pi}_j(u_i) - \pi_j(u_i))^2, \quad j = 1, \dots, C.$$

EXAMPLE 1: We generated the vectors $\mathbf{x}_i = (x_{i1}, \dots, x_{id})^\top, i = 1, \dots, n$, from a zero-mean multivariate normal with a covariance matrix $\Sigma = \{\sigma_{kl} = (0.5)^{|k-l|} : 1 \leq k, l \leq d\}$. The points u_1, \dots, u_n were generated from Uniform[0, 1]. Given each (\mathbf{x}_i, u_i) , we generated the response y_i from a Gaussian FM-VCR with $C = 2$ and the parameter setting given in Table 1, for the sample sizes $n = 200, 400$. We considered dimensions $d = 5, 10, 20, 50$, and we set non-zero coefficients $(\beta_{11}(u), \beta_{12}(u), \beta_{14}(u))$ and $(\beta_{21}(u), \beta_{23}(u))$ in the 1st and 2nd components of the mixture model. In Table 1, by increasing $d = 5$ to 10, 20, 50, the non-zero $\beta_{jl}(\cdot)$'s remain the same and we add more zeros to each component.

The average sensitivity, specificity, and estimation errors (over $R = 300$ simulated samples) for the model with $C = 2$ and $d = 5, 10, 20$ are given in Table 2. An extended version of this table which includes the results for $d = 50$ is given in Web Table 2.

[Table 2 about here.]

From Table 2, the average sensitivity and specificity corresponding to $n = 200$ vary approximately between 47% to 95% and 80% to 99%, respectively, depending on the dimension $d = 5, 10, 20$, the mixture component, and the penalty function. It seems that the task of identifying the true non-zero regression coefficients $\beta_{jl}(\cdot)$ compared to the identification of true zero coefficients, for the smaller sample size $n = 200$, is more difficult. As the sample size increases to $n = 400$, all the penalties improve in terms of identifying both true non-zero and zero regression functions. In terms of the estimation error (L_2), the SCAD and MCP are closer to the oracle estimator, followed by the AdpLASSO, LASSO, and MLLE has the worst performance as expected. Overall, the AdpLASSO, SCAD, and MCP perform similarly in terms of the three performance measures, followed by the LASSO. The results for $d = 50$ given in Web Table 2 follow the same trend with better performance for larger n .

The results for the model with $C = 3, d = 5, 10, 20, 50$ and $n = 200, 400, 600, 800$ are given in Web Tables 3 and 4. This is clearly a more challenging setting and we can see that

when $n = 200$, the sensitivity in one of the mixture components is 35% with a specificity around 80%. For $n = 800$, the sensitivity reaches 75% for this component. For the other two components, the sensitivity and specificity values are 80%-96% and 75%-92%, respectively.

Figure 1 shows the mixture order selection results based on the BIC_2 in (A.79), for the model with $C = 2$. This figure appears in color in the electronic version of this article, and any mention of color refers to that version. For a given dimension d , as n increases the performance of BIC_2 improves. As d increases from 5 to 20, the model becomes more complex and the model selection task becomes more difficult, as expected. Overall, the correct order is selected most often (and the majority is over 85%) even in the most difficult scenario.

[Figure 1 about here.]

EXAMPLE 2: Since the covariates in our real data are SNPs, we simulate x_j 's to mimic such variables. We considered discrete covariates, each taking values in $\{0, 1, 2\}$, and each was randomly generated from the multinomial distribution $\text{MNomial}((0, 1, 2), [(1-p)^2, 2p(1-p), p^2])$, where $p \in \{.05, .1, .2, .3, .4, .5\}$. We note that the smaller the value of p , the lesser the variation in the simulated covariates from this model. The parameter setting for the Gaussian FM-VCR is the same as those in Table 1. The results for $d = 5$ and $n = 200, 400$ are respectively given in Web Tables 5 and 6. The main purpose of this example is to show that for discrete covariates like SNPs, even for relatively small dimension $d = 5$, larger sample sizes are required to achieve a similar performance as those discussed in Example 1. From Web Table 5, we can see that for $n = 200$, the task of identifying non-zero regression coefficients (sensitivity) is more difficult for $p = .05$ and $.10$, which also results in higher estimation errors in one of the mixture components. As the sample size increases to $n = 400$ (Web Table 6), the overall performance of the method improves for all penalties. This result is encouraging for our real data analysis, since our sample size is even larger. In this setting, SCAD and MCP perform better than the LASSO and AdpLASSO.

7. Real data analysis

The *Fangchenggang Area Male Health and Examination Survey* was conducted to identify factors that might influence Osteocalcin (OCN) (Liao et al., 2014) and other phenotypes of interest. OCN is believed to play a role in metabolic regulation (Lee et al., 2007). It has also been observed that higher serum OCN levels are correlated with increases in bone mineral density during drug treatment for osteoporosis (Bharadwaj et al., 2009). The data include information on 2,200 unrelated healthy Chinese males, with their age, serving as our index variable (U), ranging from 20 to 69. To guard against potential confounders, individuals included in the study were free from a list of diseases, such as stroke, diabetes mellitus, primary hypertension, and hyperthyroidism. Following Zhang (2017), we focus on studying the relationship between OCN and SNPs in Chr7. We also include 7 other covariates: smoking status, physical activity, drinking, (log) body mass index, sex hormone-binding globulin, ferritin, and folic acid. After deleting individuals with missing data, SNPs with close to zero variability, and those contributing little additional information in the presence of others, we are left with $d = 43$ covariates (including 36 SNPs) and $n = 1704$ individuals.

In our analysis, we considered a Gaussian FM-VCR, where the effects of the covariates (X_j 's) on OCN (Y) and the mixing proportions (π 's) are all functions of age (U). We fitted models with $C = 1, \dots, 5$, and BIC selected $C = 2$. Since our algorithm and Example 2 based on SCAD (compared to other penalties) provided more stable estimates for different initial values of the algorithm, our results below are based on the fitted model with SCAD.

[Figure 2 about here.]

The density plots of the observed y_i 's classified into the two components of the fitted FM-VCR model are given in Figure 2(a). This figure appears in color in the electronic version of this article, and any mention of color refers to that version. Component 1 can

be interpreted as representing the high “OCN” subpopulation (hOCN), and Component 2 as the low “OCN” (lOCN) subpopulation. Figure 2(b) shows that as one ages, the probability of being in the hOCN decreases, which is consistent with biological theory since it has been found that osteocalcin is related to bone mineral density (Bharadwaj et al., 2009). From Figure 2(c,d), we observe that the estimated functional intercepts in both subpopulations are age-dependent, with decreasing OCN values as age increases; the decrease appears to slow down starting from age 35 for the hOCN, whereas the decrease is fairly linear for the lOCN. Furthermore, for the hOCN subpopulation, we found the effect of $\log(\text{BMI})$ on OCN to be negative but non-linear (varying with age), with the negative effect generally decreasing with age (Figure 2(e)). Among the other non-genetic covariates considered, folic acid also has a non-linear negative effect on OCN (Web Figures 1(d)). Most significantly, rs7456421, an SNP linked to the OCN level (Zhang, 2017), is identified to have a significant effect on OCN in the hOCN subpopulation. The effect is non-constant over time; in particular, having the SNP would negatively impact the OCN at a younger age (Figure 2(f)). Further discussion is provided in Section 8. This SNP resides in gene HIPK2, which codes for homeodomain interacting protein kinase 2, and is a multi-functional signaling molecule, including being studied as a tumor suppressor recently (Feng et al., 2017) and as a bone morphogenetic protein (Harada et al., 2003). There are a number of other SNPs that show significant effects in the hOCN subpopulation (Web Figure 1(a)-(c)). The estimated variances in both subpopulations are given in Web Figure 1(e)-(f).

In contrast, in the lOCN subpopulation, beyond the finding that the OCN level decreases with age, it is rather interesting to see that none of the other factors investigated are selected by the penalization method. Taken together, our results not only corroborated the findings of Zhang (2017) but also strengthened the results by providing evidence of a dynamic relationship in one of the two subpopulations identified by the model. In the literature, OCN

levels have been used as a biomarker for measuring osteoporosis treatment effects ([Bharadwaj et al., 2009](#)). Therefore, our finding of the differential effects of SNPs and other covariates on OCN between the two subpopulations perhaps warrants further investigation into whether the conclusion on the use of OCN as a biomarker still holds within each of these subpopulations.

8. Discussion

We developed penalized local-kernel likelihood methods for FM-VCR models and established their consistency in estimation and variable selection. We examined the finite-sample performance of the methods via simulation, and they performed well for the dimensions and sample sizes considered here. When the covariates are discrete, a larger n is required to achieve similar performance as those with continuous covariates. We observed that the method based on all four penalties performs similarly, and none of the penalties universally dominates the others. Thus, in practice, we suggest analyzing a dataset using all the penalties and choosing a final fitted sparse model that optimizes a selection criterion (e.g., BIC) and has more stable estimates based on different initial values for the EM, as done in our real data analysis.

Indeed, in our OCN data analysis, we explored several ways of analyzing the data. We fitted an FMR model which showed a lack of fit ([Web Figure 2](#)), motivating the use of an FM-VCR model. Additional details are provided in [Web Appendix I](#). In our FM-VCR analysis, we used the penalties investigated in the simulation and with multiple initial values for the EM algorithm. Our results are based on the SCAD penalty, which provided the most stable outcome and was selected by the BIC. The fit of the model to the data is seen to improve ([Web Figure 2](#)) over the FMR, and a number of age-dependent covariates were also selected. Among them, SNP rs7456421, implicated in the literature for its link to OCN level, was selected in the high OCN subpopulation and shown to have a larger impact on younger people. To further substantiate the finding of the varying effect of the SNP over the age on OCN, we added the point-wise error bars to the estimated varying coefficient plot of the SNP

(Web Figure 3(a)). We also plotted the derivative curve (approximated using Matlab) of the estimated coefficient effect over age (Web Figure 3(b)). It is evident from these plots that the underlying true derivative function is unlikely to be 0 over the entire age range, substantiating the finding of a varying coefficient for the SNP effect.

To gauge the performance of our proposed method under a number of conditions that deviate from our theoretical developments, we conducted additional simulations. First, we studied the effect of model misspecifications in terms of mixture components or the order C in (3). When the data were generated from a t FM-VCR but fitted a misspecified Gaussian FM-VCR, the results show that the performance of the method does not degrade much in variable selection and estimation error for the setting considered (Web Table 7). The details and discussion of the results are given in [Web Appendix H.1](#). On the other hand, for a model with true $C = 2$, if the order was incorrectly specified as $C = 1$, then the one-component underfitted model resembles the behavior of the larger component of the correct model, but with lower sensitivity and specificity and larger estimation errors for the corresponding measures in Component 2 when $C = 2$ was fitted (Web Table 8). In contrast, if C was incorrectly specified as larger than the true value, then the behavior of the overfitted models is similar to those with the correct model and with lower estimation errors. These simulation results are in line with the theoretical properties of over-fitted finite mixture models ([Ho et al., 2022](#)). More discussion is provided in [Web Appendix H.1](#). We recognize that our observations are based on limited simulations; thus, further study on the effect of model misspecification on our method, including properties of over-fitted FM-VCR, is warranted.

We also studied the effect of fitting an FM-VCR model when data were generated from an FMR model (without allowing for varying coefficients). The parameter setting is given in [Web Appendix H.2](#). We observe that for sensitivity, FM-VCR (the wrong model) is generally worse (on average based on 100 runs) than FMR (the true underlying model), with larger standard

deviations (SDs) (Web Table 9). On the other hand, for specificity, FM-VCR is consistently better and with much smaller SDs. This observation implies that FM-VCR selected fewer variables: fewer true positives and fewer false positives, hence larger specificity and smaller sensitivity. This observation is intuitively sensible: with a more complex model and a fixed sample size, variable selection is harder, especially for the parameter setting considered in this simulation where the effects of covariates (β_{jl} 's) are weak. Thus, in practice, one may fit both models to a data and assess the results as done in our data analysis.

In another simulation study, we also assessed the performance of the penalization method when the dimension d exceeds the sample size n , although this case is not covered in our theoretical results where d is fixed as n grows. We considered a setting with $d = 500$ and $n = 200, 400$ (Web Appendix H.4). Our results (Web Table 12) show that, while the specificity remains high and similar to the dimensions already considered in Section 6 (Table 2 and Web Table 2), there is some reduction in sensitivity, although the difference is quite small when compared to $d = 50$, indicating that reasonable results are likely to be obtained by the method even when the number of variables exceeds the sample size. Alternatively, one may screen the covariates to reduce the dimension, which appears to be reasonable as a preliminary step (Web Appendix H.4 and Web Tables 13–14). Nevertheless, rigorous study of high-dimensional settings requires new theoretical and numerical tools beyond the scope of the current paper and is a topic of future research. Finally, the mixing proportions π_j may also be considered as functions of some covariates in addition to the index variable; however, modeling and estimation become unwieldy and thus not pursued in the current study.

ACKNOWLEDGEMENTS

We would like to thank the editor, an associate editor, and two referees for their thoughtful comments. The authors also thank Professors Zengnan Mo and Ming Liao at Guangxi Medical University (Nanning, China) for their generosity of sharing the OCN data. A. Khalili and M. Asgharian are supported by the Natural Science and Engineering Research Council

of Canada (NSERC RGPIN-2020-05011) and (NSERC RGPIN-2018-05618). F. Shokoohi is supported by University of Nevada, Las Vegas, through the Startup Grant (PG18929). S. Lin is supported by NIH National Institute of General Medical Sciences (R01GM114142).

DATA AVAILABILITY

Data sharing is not applicable to this article as no new data were created or analyzed in this paper.

REFERENCES

- Berk, R., Brown, L., Buja, A., Zhang, K., and Zhao, L. (2013). Valid post-selection inference. *The Annals of Statistics* **41**, 802 – 837.
- Bharadwaj, S., Naidu, A. G. T., Betageri, G. V., Prasadaraao, N. V., and Naidu, A. S. (2009). Milk ribonuclease-enriched lactoferrin induces positive effects on bone turnover markers in postmenopausal women. *Osteoporosis International* **20**, 1603–1611.
- Dayton, C. M. and Macready, G. B. (1988). Concomitant-variable latent-class models. *Journal of the American Statistical Association* **83**, 173–178.
- Demontiero, O., Vidal, C., and Duque, G. (2012). Aging and bone loss: new insights for the clinician. *Therapeutic advances in musculoskeletal disease* **4**, 61–76.
- Dempster, A., Laird, N., and Rubin, D. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B* **39**, 1–38.
- Fan, J. and Gigbels, I. (1996). *Local Polynomial Modelling and Its Applications*. Chapman & Hall/CRC Monographs on Statistics and Applied Probability.
- Fan, J. and Li, R. (2001). Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties. *Journal of the American Statistical Association* **96**, 1348–1360.
- Fan, J. and Zhang, W. (1999). Statistical estimation in varying coefficient models. *The Annals of Statistics* **27**, 1491–1518.
- Feng, Y., Zhou, L., Sun, X., and Li, Q. (2017). Homeodomain-interacting protein kinase 2 (hipk2): a promising target for anti-cancer therapies. *Oncotarget* **8**, 20452–20461.

- Frühwirth-Schnatter, S., Celeux, G., and Robert, C. P. (2018). *Handbook of Mixture Analysis*. Chapman and Hall/CRC, 1st edition.
- Harada, J., Kokura, K., Kanei-Ishii, C., Nomura, T., Khan, M. M., Kim, Y., et al. (2003). Requirement of the co-repressor homeodomain-interacting protein kinase 2 for ski-mediated inhibition of bone morphogenetic protein-induced transcriptional activation. *Journal of Biological Chemistry* **278**, 38998–39005.
- Hastie, T. and Tibshirani, R. (1993). Varying-coefficient models. *Journal of the Royal Statistical Society. Series B (Methodological)* **55**, 757–796.
- Ho, N., Yang, C.-Y., and Jordan, M. I. (2022). Convergence rates for gaussian mixtures of experts. *Journal of Machine Learning Research* .
- Huang, M., Li, R., and Wang, S. (2013). Nonparametric mixture of regression models. *Journal of the American Statistical Association* **108**, 929–941.
- Huang, M., Yao, W., Wang, S., and Chen, Y. (2018). Statistical inference and applications of mixture of varying coefficient models. *Scandinavian Journal of Statistics* **45**, 618–643.
- Janson, S. (1987). Maximal spacing in several dimensions. *Annals of Probability* **15**, 274–80.
- Khalili, A. and Chen, J. (2007). Variable Selection in Finite Mixture of Regression Models. *Journal of the American Statistical Association* **102**, 1025–1038.
- Lee, N. K., Sowa, H., Hinoi, E., Ferron, M., Ahn, J. D., Confavreux, C., et al. (2007). Endocrine regulation of energy metabolism by the skeleton. *Cell* **130**, 456–469.
- Liao, M., Shi, J., Huang, L., Gao, Y., Tan, A., Wu, C., et al. (2014). Genome-wide association study identifies variants in pms1 associated with serum ferritin in a chinese population. *PloS one* **9**, e105844.
- Liu, D. M., Guo, X. Z., Tong, H. J., Tao, B., Sun, L.-H., Zhao, H.-Y., et al. (2015). Association between osteocalcin and glucose metabolism: a meta-analysis. *Osteoporosis International* **26**, 2823–2833.

- McLachlan, G. J. and Peel, D. (2000). *Finite mixture models*. Wiley & Sons, New York.
- Shokoohi, F., Khalili, A., Asgharian, M., and Lin, S. (2019). Capturing heterogeneity of covariate effects in hidden subpopulations in the presence of censoring and large number of covariates. *Annals of Applied Statistics* **13**, 444–465.
- Silverman, B. W. (1984). Spline Smoothing: The Equivalent Variable Kernel Method. *The Annals of Statistics* **12**, 898 – 916.
- Städler, N., Bühlmann, P., and van de Geer, S. (2010). L_1 -penalization for Mixture Regression Models. *Test* **19**, 209–256.
- Tibshirani, R. (1996). Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* **58**, 267–288.
- Wang, H. and Xia, Y. (2009). Shrinkage estimation of the varying coefficient model. *Journal of the American Statistical Association* **104**, 747–757.
- Wei, F., Huang, J., and Li, H. (2011). Variable selection and estimation in high-dimensional varying-coefficient models. *Statistica Sinica* **21**, 1515–1540.
- Xiang, S., Yao, W., and Yang, G. (2019). An overview of semiparametric extensions of finite mixture models. *Statistical Science* **34**, 391–404.
- Yuan, M. and Lin, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B* **68**, 49–67.
- Zhang, C.-H. (2010). Nearly Unbiased Variable Selection under Minimax Concave Penalty. *The Annals of Statistics* **38**, 894–942.
- Zhang, H. (2017). *Detecting rare haplotype-environmental interaction and nonlinear effects of rare haplotypes using Bayesian lasso on quantitative traits*. PhD thesis, Ohio State U.
- Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* **101**, 1418–1429.

SUPPORTING INFORMATION

Web Appendices, Tables, and Figures referenced in Sections 2-8, and codes used for simulation and data analysis, are available with this paper at the Biometrics website on Wiley Online Library.

Received January 2022. Revised August 2022. Accepted September 2022.

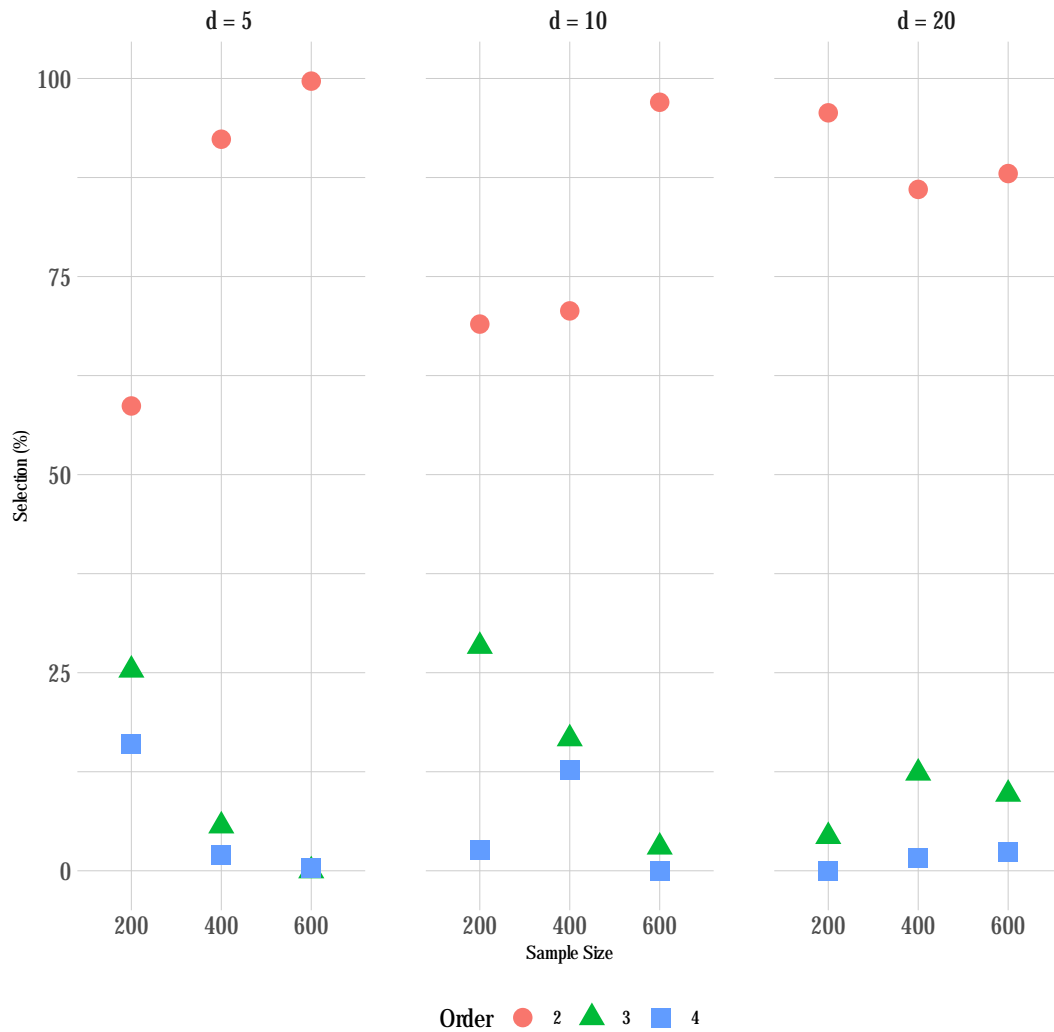


Figure 1: Example 1: Order selection results for the model with true order $C = 2$. This figure appears in color in the electronic version of this article, and any mention of color refers to that version.

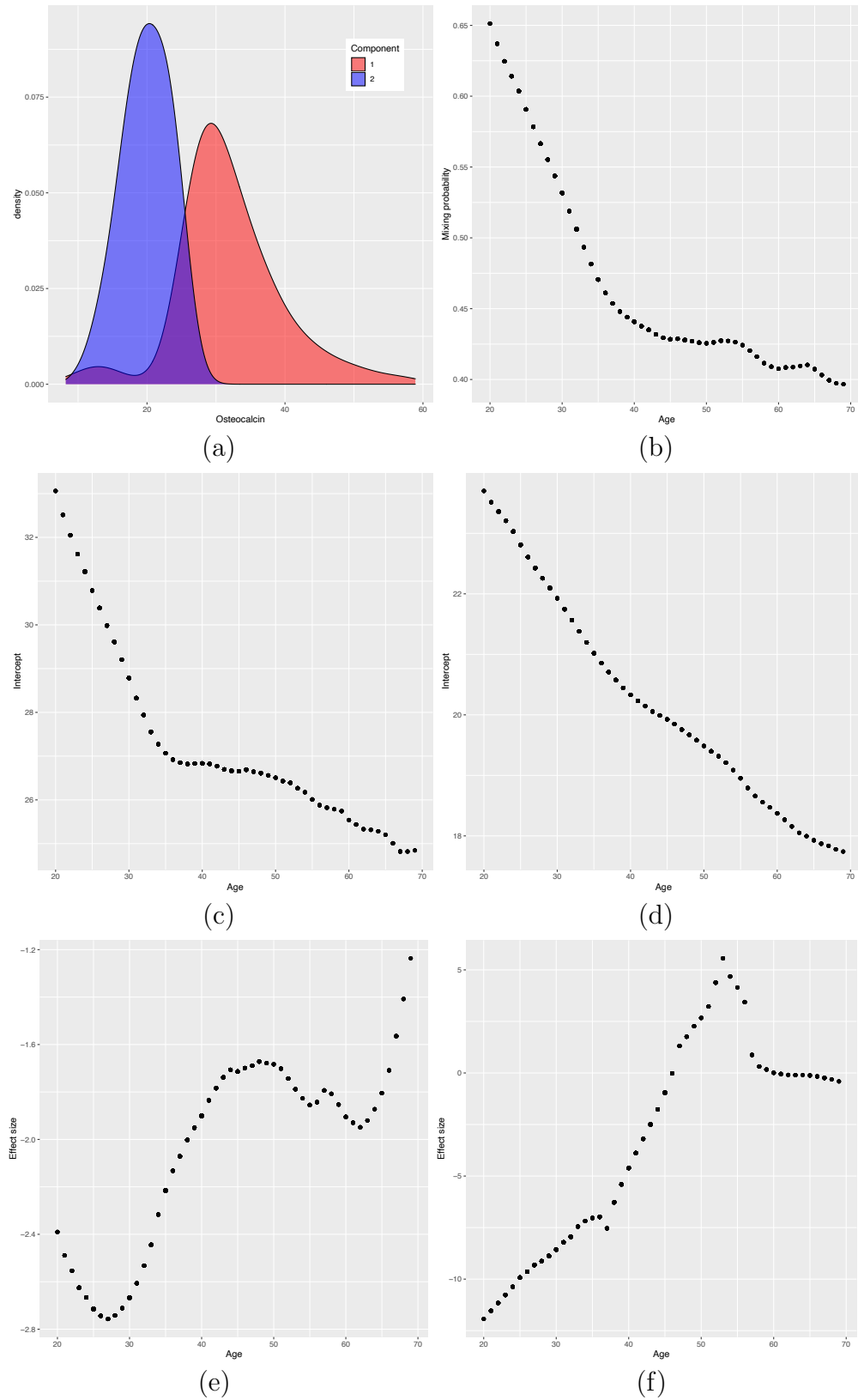


Figure 2: Osteocalcin data analysis; (a) Density plots: Component 1 (hOCN, red), Component 2 (lOCN, blue); (b) Estimated mixing probabilities over time in hOCN; (c) Estimated intercept over time in hOCN; (d) Estimated intercept over time in lOCN; (e) Estimated effect of log(BMI) over time in hOCN; (f) Estimated effect of rs7456421 over time in hOCN. This figure appears in color in the electronic version of this article, and any mention of color refers to that version.

Table 1: Parameters settings for the Gaussian FM-VCR model with $C = 2$.

Component(j):	1	2
parameters	$d(= 5, 10, 20, 50)^1$	
$\beta_{j0}(u)$	-2	-1
$\beta_{j1}(u)$	$1 + 0.5 \cos(\pi u)$	$1.5 \sin(\pi u)$
$\beta_{j2}(u)$	$1 + 0.5 \cos(2\pi u)$	0
$\beta_{j3}(u)$	0	$1.5 - 0.5 \sin(\pi u/2)$
$\beta_{j4}(u)$	$\sin(6\pi u)$	0
$\beta_{j5}(u)$	0	0
\vdots	\vdots	\vdots
$\beta_{j,50}(u)$	0	0
$\sigma_j(u)$	$0.3e^{(0.5u)}$	$0.5e^{(-0.2u)}$
$\pi_j(u)$	$e^{0.5u}/(1 + e^{0.5u})$	$(1 + e^{0.5u})^{-1}$

¹ We keep the non-zero coefficient functions $\beta_{j\ell}(\cdot)$ the same for different values of the dimension d .

Table 2: Results of Example 1: average (SD) sensitivity, specificity, and estimation errors.

$C = 2$		Criteria Component	Sensitivity		Specificity		$L_2(\hat{\beta}_j)$		$L_2(\hat{\sigma}_j^2)$		$L_2(\hat{\pi}_j)$	
d	n		1	2	1	2	1	2	1	2	1	2
5	200	Oracle	—	—	—	—	.650	.162	.383	.104	.321	.104
		MLE	—	—	—	—	1.04	.243	.748	.127	.370	.127
		AdpLASSO	.687(.203)	.948(.191)	.967(.125)	.962(.122)	.934	.208	.612	.118	.500	.118
		LASSO	.701(.223)	.847(.347)	.813(.265)	.824(.214)	.983	.256	.812	.167	.524	.167
		MCP	.682(.127)	.952(.153)	.978(.110)	.972(.120)	.888	.163	.561	.094	.505	.094
	SCAD	.684(.127)	.957(.147)	.980(.106)	.974(.118)	.882	.165	.554	.093	.507	.093	
	400	Oracle	—	—	—	—	.399	.083	.214	.080	.273	.080
		MLE	—	—	—	—	.547	.097	.389	.075	.278	.075
		AdpLASSO	.978(.083)	.997(.058)	.998(.030)	.997(.043)	.499	.082	.225	.074	.311	.074
		LASSO	.989(.060)	.995(.064)	.938(.179)	.910(.180)	.582	.116	.376	.086	.338	.086
MCP		.983(.073)	1.00(.000)	1.00(.000)	.999(.019)	.458	.078	.208	.073	.313	.073	
SCAD	.984(.070)	1.00(.000)	1.00(.000)	.999(.019)	.458	.077	.207	.073	.313	.073		
10	200	Oracle	—	—	—	—	.698	.162	.397	.089	.314	.089
		MLE	—	—	—	—	1.23	.287	1.02	.138	.335	.138
		AdpLASSO	.589(.239)	.893(.259)	.970(.067)	.966(.061)	1.09	.288	.848	.161	.501	.161
		LASSO	.470(.321)	.775(.401)	.928(.100)	.917(.087)	1.32	.383	1.15	.294	.505	.294
		MCP	.581(.194)	.895(.245)	.984(.045)	.981(.048)	1.04	.238	.763	.106	.503	.106
	SCAD	.594(.188)	.892(.243)	.983(.049)	.980(.050)	1.03	.237	.761	.107	.501	.107	
	400	Oracle	—	—	—	—	.439	.089	.227	.082	.273	.082
		MLE	—	—	—	—	.878	.149	.607	.087	.337	.087
		AdpLASSO	.690(.127)	.993(.057)	.997(.023)	.995(.029)	.808	.165	.390	.069	.411	.069
		LASSO	.766(.230)	.970(.166)	.970(.066)	.962(.064)	.862	.232	.518	.108	.425	.108
MCP		.692(.104)	.992(.064)	.997(.020)	.998(.014)	.789	.145	.374	.06	.412	.062	
SCAD	.688(.098)	.992(.064)	.996(.023)	.998(.016)	.792	.146	.377	.062	.412	.062		
20	200	Oracle	—	—	—	—	.615	.095	.297	.079	.313	.079
		MLE	—	—	—	—	1.37	.225	1.30	.082	.383	.082
		AdpLASSO	.646(.115)	.943(.202)	.981(.043)	.984(.042)	.914	.269	.713	.056	.488	.056
		LASSO	.657(.120)	.953(.187)	.963(.059)	.952(.053)	.924	.336	.896	.056	.494	.056
		MCP	.643(.121)	.953(.199)	.994(.028)	.995(.022)	.915	.220	.619	.052	.467	.052
	SCAD	.641(.127)	.947(.206)	.994(.028)	.996(.021)	.918	.227	.633	.052	.472	.052	
	400	Oracle	—	—	—	—	.353	.072	.211	.067	.251	.067
		MLE	—	—	—	—	.912	.102	.616	.038	.423	.038
		AdpLASSO	.668(.019)	.998(.038)	.999(.015)	1.00(.005)	.798	.170	.417	.036	.484	.036
		LASSO	.716(.118)	1.00(.000)	.978(.035)	.969(.041)	.823	.237	.523	.036	.499	.036
MCP		.667(.047)	.997(.058)	.999(.010)	1.00(.007)	.796	.159	.433	.036	.481	.036	
SCAD	.668(.051)	.997(.058)	.999(.011)	1.00(.007)	.796	.159	.433	.036	.481	.036		

**Supporting Information for “Sparse Estimation in Semi-parametric Finite
Mixture of Varying Coefficient Regression Models” by**

Abbas Khalili^{1,*}, Farhad Shokoohi^{2,}, Masoud Asgharian^{3,***}, and Shili Lin^{4,****}**

^{1,3}Department of Mathematics and Statistics, McGill University, Montreal, QC, Canada

²Department of Mathematical Sciences, University of Nevada Las Vegas, Las Vegas, NV 89154, USA

⁴Department of Statistics, Ohio State University, Columbus, OH 43210, USA

**email:* abbas.khalili@mcgill.ca

***email:* farhad.shokoohi@unlv.edu

****email:* masoud.asgharian2@mcgill.ca

*****email:* shili@stat.osu.edu

Web Appendix A. Identifiability conditions of f_C^*

The following conditions are considered for an FM-VCR model with probability density (mass) function f_C^* in (3) of the paper to be identifiable (see Theorem 1 of Huang et al. (2018)).

(C1) For $j = 1, \dots, C$, $\pi_j(u) > 0$ is a continuous function, and $\beta_j(u)$ and $\phi_j(u)$ are differentiable functions of u .

(C2) Any two curves $\mathbf{a}_1(u) = (\beta_i(u), \phi_i(u))$ and $\mathbf{a}_2(u) = (\beta_j(u), \phi_j(u))$, $i \neq j$, are transversal; i.e., for any $u \in \mathcal{U}$, $\|\mathbf{a}_1(u) - \mathbf{a}_2(u)\|^2 + \|\mathbf{a}'_1(u) - \mathbf{a}'_2(u)\|^2 \neq 0$, and $\|\cdot\|$ is the Euclidean norm.

(C3) The support for U , denoted by \mathcal{U} , is a compact subset of the real numbers \mathbb{R} .

(C4) The density $g(\mathbf{x}|u)$ of \mathbf{X} given $U = u$ is a full dimensional (d -dimensional) density function.

Web Appendix B. Regularity Conditions

In the following three subsections, we state Regularity Conditions **(RC.1)-(RC.5)** on the distributions f_C^* , g , and m which is the density of U ; Conditions **(KC.1)-(KC.2)** on the kernel K , and Conditions **(PC.1)-(PC.3)** on the penalty p_n and the smoothing/tuning parameters h_n and λ_n .

Web Appendix B.1 Conditions on f_C^* , g , and m

(RC.1) The set $\{(U_i, \mathbf{X}_i, Y_i), i = 1, \dots, n\}$ is a sample of independent and identically distributed variables drawn from its population $\mathbf{V} = (U, \mathbf{X}, Y)$ with the probability density function $f_{\mathbf{V}}(u, \mathbf{x}, y) = f_C^*(y|\psi(u), \mathbf{x})g(\mathbf{x}|u)m(u)$. The density has a common support in \mathbf{v} for all the parameter values $\psi \in \Omega$, and $f_{\mathbf{V}}$ or equivalently the mixture density f_C^* in (3) of the paper is identifiable up to a permutation of the mixture components. (See Section **Web Appendix A** for identifiability.)

(RC.2) The density $m(u)$ of U is positively bounded away from zero over its compact support \mathcal{U} ; it is continuously twice differentiable, and these derivatives are uniformly bounded over \mathcal{U} . The density f_C^* admits third partial derivatives with respect to its parameters $\boldsymbol{\psi}$ for all (u, \boldsymbol{x}, y) ; g admits its third derivatives with respect to u , for all \boldsymbol{x} . The function $\boldsymbol{\psi}(u)$ has continuous third derivatives over \mathcal{U} .

(RC.3) There exist functions $M_k(\boldsymbol{x}, y)$, $k = 1, 2$, with $E_{(\boldsymbol{X}, Y)}\{M_1^{m_0}(\boldsymbol{X}, Y)\} < \infty$ for some $m_0 > 2$, and $E_{(\boldsymbol{X}, Y)}\{M_2(\boldsymbol{X}, Y)\} < \infty$, and $\int_{\boldsymbol{x}, y} M_1(\boldsymbol{x}, y) d\boldsymbol{x} dy < \infty$, such that for all \boldsymbol{x} and y , and all $\boldsymbol{\psi}(u)$ in a neighbourhood of $\boldsymbol{\psi}^0(u)$,

$$\left| \frac{\partial l(\boldsymbol{\psi}, \boldsymbol{x}, y)}{\partial \psi_j} \right| < M_1(\boldsymbol{x}, y), \quad \left| \frac{\partial^2 l(\boldsymbol{\psi}, \boldsymbol{x}, y)}{\partial \psi_j \partial \psi_k} \right| < M_1(\boldsymbol{x}, y),$$

and

$$\left| \frac{\partial^3 l(\boldsymbol{\psi}, \boldsymbol{x}, y)}{\partial \psi_j \partial \psi_k \partial \psi_l} \right| < M_2(\boldsymbol{x}, y).$$

(RC.4) The matrix $\boldsymbol{I}(u)$ in (12) of the paper is continuous in u and positive definite for all $u \in \mathcal{U}$; the vector $\Delta(u; u_t)$ in (13) of the paper is continuously third time differentiable with respect to u over \mathcal{U} .

(RC.5) The conditions

$$E_0 \left(\left| \frac{\partial l(\boldsymbol{\psi}, \boldsymbol{X}, Y)}{\partial \psi_j} \right|^3 \right) < \infty \quad \text{and} \quad E_0 \left(\left| \frac{\partial^2 l(\boldsymbol{\psi}, \boldsymbol{X}, Y)}{\partial \psi_i \partial \psi_j} \right|^{2+\delta} \right) < \infty$$

hold for all i and j and some $\delta > 0$, where E_0 is the expectation with respect to the joint distribution of (Y, \boldsymbol{X}, U) .

Web Appendix B.2 *Conditions on the kernel K*

(KC.1) K is a Lipschitz continuous and symmetric pdf with a compact support \mathcal{U} .

(KC.2) $\mathcal{K}_i = \int_{\mathcal{U}} t^i K(t) dt$, $\mathcal{V}_i = \int_{\mathcal{U}} t^i K^2(t) dt$, $i = 0, \dots, 4$, are finite; for i odd they are zero.

Web Appendix B.3 Conditions on the penalty p_n , and the tuning parameters (λ_n, h_n)

With regard to the penalty function p_n , we denote the quantities

$$r_{1n} = \max_{\substack{1 \leq l \leq d_j^0 \\ 1 \leq j \leq C}} \left\{ \frac{h_n^{1/2} |p'_n(\theta_{jl}^0; \lambda_n)|}{n^{3/2}} : \theta_{jl}^0 = \sqrt{E_0\{\beta_{jl}^0(U_t)\}^2} \neq 0 \right\}, \text{ and}$$

$$r_{2n} = \max_{\substack{1 \leq l \leq d_j^0 \\ 1 \leq j \leq C}} \left\{ \frac{|p''_n(\theta_{jl}^0; \lambda_n)|}{n^2} : \theta_{jl}^0 = \sqrt{E_0\{\beta_{jl}^0(U_t)\}^2} \neq 0 \right\},$$

where $p'_n(\theta_{jl}^0; \lambda_n)$ and $p''_n(\theta_{jl}^0; \lambda_n)$ are the first and second derivatives of $p_n(\theta_{jl}; \lambda_n)$ evaluated at $\theta_{jl} = \theta_{jl}^0 \neq 0$.

(PC.1) $p_n(\theta; \lambda_n)$ is nonnegative and symmetric in $\theta \in \mathbb{R}$, and $p_n(0; \lambda_n) = 0$. It is also nondecreasing and twice continuously differentiable for all but finitely many values $\theta \in (0, \infty)$.

(PC.2) As $n \rightarrow \infty$, $\lambda_n = o(1)$ such that $\min_{j,l} \theta_{jl}^0 / \lambda_n \rightarrow \infty$; and $h_n = o(1)$ such that $nh_n \rightarrow \infty$, and $nh_n^5 = O(1)$. Also, $r_{1n} = o(\sqrt{nh_n})$ and $r_{2n} = o(1)$.

(PC.3) Let $N_{h_n} = \left\{ \theta; 0 < |\theta| \leq \frac{\log(nh_n)}{\sqrt{nh_n}} \right\}$. Then, $\liminf_{n \rightarrow \infty} \inf_{\theta \in N_{h_n}} h_n^{1/2} \frac{p'_n(\theta; \lambda_n)}{n^{3/2}} = \infty$.

The smoothness Condition **(RC.1)** facilitates obtaining estimators of $\beta_{jl}(\cdot)$ by differentiating $\tilde{L}_n(\Psi; \lambda, h)$ and for studying the asymptotic properties of the estimators of the true non-zero $\beta_{jl}(\cdot)$. **(PC.2)** is required to obtain $\sqrt{nh_n}$ -consistent estimators of the true non-zero $\beta_{jl}(\cdot)$, while **(PC.3)** is required for sparsity.

Web Appendix C. Examples of the penalty function p_n and the kernel K

We used the following penalties and kernel function in our simulations and analysis of the real data.

- LASSO: $p_n(\theta; \lambda) / n^2 = \lambda |\theta|$;
- AdpLASSO: $p_n(\theta; \lambda) / n^2 = \lambda w |\theta|$, for some (possibly random) known weights w ;
- MCP: $p'_n(\theta; \lambda) / n^2 = \text{sgn}(\theta) \frac{(a\lambda - |\theta|)_+}{a}$;

- SCAD: $p'_n(\theta; \lambda)/n^2 = \text{sgn}(\theta) \left\{ \lambda I(|\theta| \leq \lambda) + \frac{(a\lambda - |\theta|)_+}{a-1} I(|\theta| > \lambda) \right\},$

where $p'_n(\cdot; \lambda)$ is the first derivative of a penalty with respect to θ , and $(x)_+ = \max\{0, x\}$.

The Epanechnikov Kernel: $K(t) = \frac{3}{4}(1 - u^2)_+.$

Web Appendix D. Proofs

The result of the following lemma is used in the sequel.

LEMMA 1: (**Barlett's first and second identities**) *Under Conditions (RC.1)-(RC.4) in Web Appendix B, and for any interior point u_t of \mathcal{U} , it holds that*

$$E_{(Y, \mathbf{X})|U} [l'(\boldsymbol{\psi}^0(U), \mathbf{X}, Y)|U = u_t] = 0,$$

and

$$\begin{aligned} & E_{(Y, \mathbf{X})|U} [l''(\boldsymbol{\psi}^0(U), \mathbf{X}, Y)|U = u_t] \\ &= -E_{(Y, \mathbf{X})|U} [l'(\boldsymbol{\psi}^0(U), \mathbf{X}, Y)l'^\top(\boldsymbol{\psi}^0(U), \mathbf{X}, Y)|U = u_t] = \mathbf{I}(u_t). \end{aligned}$$

Proof. Conditioning on $U = u_t$, we prove Barlett's first identity as follows.

$$\begin{aligned} & E_{(Y, \mathbf{X})|U} [l'(\boldsymbol{\psi}^0(U), \mathbf{X}, Y)|U = u_t] \\ &= \int_{\mathcal{Y}} \int_{\mathcal{X}} l'(\boldsymbol{\psi}^0(u_t), \mathbf{x}, y) f_C^*(y|\boldsymbol{\psi}^0(u_t), \mathbf{x}) g(\mathbf{x}|u_t) d\mathbf{x} dy \\ &= \int_{\mathcal{Y}} \int_{\mathcal{X}} \frac{\partial l(\boldsymbol{\psi}, \mathbf{x}, y)}{\partial \boldsymbol{\psi}} \Big|_{\boldsymbol{\psi}=\boldsymbol{\psi}^0(u_t)} f_C^*(y|\boldsymbol{\psi}^0(u_t), \mathbf{x}) g(\mathbf{x}|u_t) d\mathbf{x} dy \\ &= \int_{\mathcal{Y}} \int_{\mathcal{X}} \frac{\partial \log f_C^*(y|\boldsymbol{\psi}, \mathbf{x})}{\partial \boldsymbol{\psi}} \Big|_{\boldsymbol{\psi}=\boldsymbol{\psi}^0(u_t)} f_C^*(y|\boldsymbol{\psi}^0(u_t), \mathbf{x}) g(\mathbf{x}|u_t) d\mathbf{x} dy \\ &= \int_{\mathcal{Y}} \int_{\mathcal{X}} \frac{\partial f_C^*(y|\boldsymbol{\psi}, \mathbf{x})}{\partial \boldsymbol{\psi}} \Big|_{\boldsymbol{\psi}=\boldsymbol{\psi}^0(u_t)} g(\mathbf{x}|u_t) d\mathbf{x} dy \\ &= \frac{\partial}{\partial \boldsymbol{\psi}} \int_{\mathcal{Y}} \int_{\mathcal{X}} \left[\sum_{c=1}^C \pi_c f(y; \theta_c(\mathbf{x}), \phi_c) \right] g(\mathbf{x}|u_t) d\mathbf{x} dy \Big|_{\boldsymbol{\psi}=\boldsymbol{\psi}^0(u_t)} \\ &= \frac{\partial}{\partial \boldsymbol{\psi}} (1) = 0. \end{aligned}$$

For a proof of Barlett’s second identity note that

$$\begin{aligned}
& - E_{(Y, \mathbf{X})|U} [l''(\boldsymbol{\psi}^0(U), \mathbf{X}, Y)|U = u_t] \\
&= - \int_{\mathcal{Y}} \int_{\mathcal{X}} \frac{\partial^2 \log f_C^*(y|\boldsymbol{\psi}, \mathbf{x})}{\partial \boldsymbol{\psi} \partial \boldsymbol{\psi}^\top} \Big|_{\boldsymbol{\psi}=\boldsymbol{\psi}^0(u_t)} f_C^*(y|\boldsymbol{\psi}^0(u_t), \mathbf{x}) g(\mathbf{x}|u_t) d\mathbf{x} dy \\
&= - \int_{\mathcal{Y}} \int_{\mathcal{X}} \frac{\frac{\partial^2 f_C^*(y|\boldsymbol{\psi}, \mathbf{x})}{\partial \boldsymbol{\psi} \partial \boldsymbol{\psi}^\top} f_C^*(y|\boldsymbol{\psi}, \mathbf{x}) - \left\{ \frac{\partial f_C^*(y|\boldsymbol{\psi}, \mathbf{x})}{\partial \boldsymbol{\psi}} \right\} \left\{ \frac{\partial f_C^*(y|\boldsymbol{\psi}, \mathbf{x})}{\partial \boldsymbol{\psi}} \right\}^\top}{\{f_C^*(y|\boldsymbol{\psi}, \mathbf{x})\}^2} \Big|_{\boldsymbol{\psi}=\boldsymbol{\psi}^0(u_t)} \\
&\quad \times f_C^*(y|\boldsymbol{\psi}^0(u_t), \mathbf{x}) g(\mathbf{x}|u_t) d\mathbf{x} dy \\
&= - \int_{\mathcal{Y}} \int_{\mathcal{X}} \frac{\partial^2}{\partial \boldsymbol{\psi} \partial \boldsymbol{\psi}^\top} \{f_C^*(y|\boldsymbol{\psi}, \mathbf{x})\} \Big|_{\boldsymbol{\psi}=\boldsymbol{\psi}^0(u_t)} g(\mathbf{x}|u_t) d\mathbf{x} dy \\
&\quad + \int_{\mathcal{Y}} \int_{\mathcal{X}} \frac{\left\{ \frac{\partial}{\partial \boldsymbol{\psi}} f_C^*(y|\boldsymbol{\psi}, \mathbf{x}) \right\} \left\{ \frac{\partial}{\partial \boldsymbol{\psi}} f_C^*(y|\boldsymbol{\psi}, \mathbf{x}) \right\}^\top}{\{f_C^*(y|\boldsymbol{\psi}, \mathbf{x})\}^2} \Big|_{\boldsymbol{\psi}=\boldsymbol{\psi}^0(u_t)} f_C^*(y|\boldsymbol{\psi}^0(u_t), \mathbf{x}) g(\mathbf{x}|u_t) d\mathbf{x} dy \\
&= - \frac{\partial^2}{\partial \boldsymbol{\psi} \partial \boldsymbol{\psi}^\top} \int_{\mathcal{Y}} \int_{\mathcal{X}} \{f_C^*(y|\boldsymbol{\psi}, \mathbf{x})\} \Big|_{\boldsymbol{\psi}=\boldsymbol{\psi}^0(u_t)} g(\mathbf{x}|u_t) d\mathbf{x} dy \\
&\quad + \int_{\mathcal{Y}} \int_{\mathcal{X}} l'(\boldsymbol{\psi}^0(u_t), \mathbf{x}, y) l'^\top(\boldsymbol{\psi}^0(u_t), \mathbf{x}, y) f_C^*(y|\boldsymbol{\psi}^0(u_t), \mathbf{x}) g(\mathbf{x}|u_t) d\mathbf{x} dy \\
&= E_{(Y, \mathbf{X})|U} \{l'(\boldsymbol{\psi}^0(u_t), \mathbf{X}, Y) l'^\top(\boldsymbol{\psi}^0(u_t), \mathbf{X}, Y)|U = u_t\}.
\end{aligned}$$

We now provide the proofs of our main results, starting with a proposition on point-wise consistency of the maximum local-kernel log-likelihood estimator (MLLE) of the parameter vector of $\boldsymbol{\psi}(u)$ in the FM-VCR model in (3) of the paper .

PROPOSITION 1: Suppose that Regularity Conditions (RC.1)-(RC.4) on the parametric family and Conditions (KC.1)-(KC.2) on the kernel K stated in Web Appendix B of this document hold, and $h_n = o(1)$ such that $nh_n \rightarrow \infty$, and $nh_n^5 = O(1)$. Then for any $u \in \mathcal{U}$, there exists a local maximizer $\check{\boldsymbol{\psi}}_n(u)$ of $\ell_n(\boldsymbol{\psi}(u); h_n)$ in (7) of the paper such that: $\|\check{\boldsymbol{\psi}}_n(u) - \boldsymbol{\psi}^0(u)\|_2 = O_p\{(nh_n)^{-1/2}\}$.

Proof of Proposition 1

Proof. The proof is a special case of the proof of Theorem 1 (without the penalty) and thus omitted.

Proof of Theorem 1

Proof. First note that, for any $u \in \mathcal{U}$,

$$\boldsymbol{\psi}(u) = (\boldsymbol{\pi}^\top(u), \boldsymbol{\phi}^\top(u), \boldsymbol{\beta}_1^\top(u), \dots, \boldsymbol{\beta}_C^\top(u))^\top \in \mathbb{R}^{(C(d+2)-1) \times 1},$$

and whenever necessary, we rewrite

$$\boldsymbol{\psi}(u) = (\psi_1(u), \dots, \psi_{C(d+2)-1}(u))^\top$$

without changing the order of $\boldsymbol{\phi}^\top(u), \boldsymbol{\pi}^\top(u), \boldsymbol{\beta}_1^\top(u), \dots, \boldsymbol{\beta}_C^\top(u)$. Otherwise, we will use the same notation as defined in Section 2 of the paper. Let

$$\boldsymbol{w}_{j,I} = \begin{pmatrix} w_{j1,I,1} & \dots & w_{jd,I,1} \\ \vdots & \ddots & \vdots \\ w_{j1,I,n} & \dots & w_{jd,I,n} \end{pmatrix} \in \mathbb{R}^{n \times d},$$

and the vectors

$$\boldsymbol{w}_{j,II} = (w_{j,II,1}, \dots, w_{j,II,n})^\top \in \mathbb{R}^{n \times 1}, \quad j = 1, \dots, C,$$

$$\boldsymbol{w}_{j,III} = (w_{j,III,1}, \dots, w_{j,III,n})^\top \in \mathbb{R}^{n \times 1}, \quad j = 1, \dots, C-1.$$

We set

$$\boldsymbol{W}_n = (\boldsymbol{w}_{1,I}, \dots, \boldsymbol{w}_{C,I}, \boldsymbol{w}_{1,II}, \dots, \boldsymbol{w}_{C,II}, \boldsymbol{w}_{1,III}, \dots, \boldsymbol{w}_{C-1,III}) \in \mathbb{R}^{n \times (C(d+2)-1)}$$

which is an arbitrary matrix with rows \boldsymbol{w}_t^\top , and $w_{j,I,t}$ is denoted as the t -th row of $\boldsymbol{w}_{j,I}$, for $t = 1, \dots, n$. The first Cd columns of \boldsymbol{w} are denoted by $\boldsymbol{v}_{11}, \dots, \boldsymbol{v}_{1d}, \dots, \boldsymbol{v}_{C1}, \dots, \boldsymbol{v}_{Cd}$. For an arbitrary matrix $\boldsymbol{A} = (a_{ij})$, we define its L_2 norm (i.e. the Frobenius norm) as

$$\|\boldsymbol{A}\|_F = \left(\sum a_{ij}^2 \right)^{1/2}.$$

Note that $\|\mathbf{W}_n\|_F^2 = \sum_{t=1}^n \|\mathbf{w}_t\|^2$ and $\sum_{j=1}^C \sum_{l=1}^d \|\mathbf{v}_{jl}\|^2 \leq \|\mathbf{W}_n\|_F^2$. Whenever necessary, we rewrite $\mathbf{w}_t = (w_{1,t}, \dots, w_{C(d+2)-1,t})^\top$; otherwise, we will use the same notation as defined by

$$\mathbf{w}_t = (w_{11,I,t}, \dots, w_{1d,I,t}, \dots, w_{C1,I,t}, \dots, w_{Cd,I,t}, w_{1,II,t}, \dots, w_{C,II,t}, w_{1,III,t}, \dots, w_{C-1,III,t})^\top.$$

To prove the claim of the theorem, it suffices to show that for any small $\epsilon > 0$, there exists a constant M_ϵ such that

$$\lim_{n \rightarrow \infty} P \left\{ \sup_{n^{-1}\|\mathbf{W}_n\|_2^2 = M_\epsilon^2} \tilde{L}_n(\Psi_n^0 + \gamma_n \mathbf{W}_n; \boldsymbol{\lambda}_n, h_n) < \tilde{L}_n(\Psi_n^0; \boldsymbol{\lambda}_n, h_n) \right\} \geq 1 - \epsilon, \quad (\text{A.1})$$

where $\gamma_n = (1 + r_{1n})(nh_n)^{-1/2}$.

To show (A.1), we proceed as follows. Let

$$D_n(\mathbf{W}_n) = h_n n^{-1} \left\{ \tilde{L}_n(\Psi_n^0 + \gamma_n \mathbf{W}_n; \boldsymbol{\lambda}_n, h_n) - \tilde{L}_n(\Psi_n^0; \boldsymbol{\lambda}_n, h_n) \right\}.$$

By the definition of penalized local log-likelihood \tilde{L}_n in (9) of the paper, we have

$$\tilde{L}_n(\Psi_n^0 + \gamma_n \mathbf{W}_n; \boldsymbol{\lambda}_n, h_n) = \sum_{t=1}^n \sum_{i=1}^n l(\boldsymbol{\psi}^0(U_t) + \gamma_n \mathbf{w}_t, \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t) - \mathbb{P}_n(\Psi_n^0 + \gamma_n \mathbf{W}_n; \boldsymbol{\lambda}_n)$$

and

$$\tilde{L}_n(\Psi_n^0; \boldsymbol{\lambda}_n, h_n) = \sum_{t=1}^n \sum_{i=1}^n l(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t) - \mathbb{P}_n(\Psi_n^0; \boldsymbol{\lambda}_n),$$

where $l(\boldsymbol{\psi}^0(u), \mathbf{x}, y) = \log \{f_C^*(y|\boldsymbol{\psi}^0(u), \mathbf{x})\}$, f_C^* is the mixture density given in (3) of the paper, and the penalty

$$\mathbb{P}_n(\Psi_n; \boldsymbol{\lambda}_n) = \sum_{j=1}^C \sum_{l=1}^d p_n(\|\mathbf{b}_{jl}\|_2 / \sqrt{n}; \lambda_{nj}).$$

By Condition (PC.1) on the penalty, we have $p_n(0; \lambda_{nj}) = 0$. Thus,

$$\begin{aligned} D_n(\mathbf{W}_n) &= h_n n^{-1} \sum_{t=1}^n \sum_{i=1}^n \left\{ l(\boldsymbol{\psi}^0(U_t) + \gamma_n \mathbf{w}_t, \mathbf{X}_i, Y_i) - l(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) \right\} K_{h_n}(U_i - U_t) \\ &\quad - h_n n^{-1} \sum_{j=1}^C \sum_{l=1}^d \left\{ p_n(\|\mathbf{b}_{jl}^0 + \gamma_n \mathbf{v}_{jl}\|_2 / \sqrt{n}; \lambda_{nj}) - p_n(\|\mathbf{b}_{jl}^0\|_2 / \sqrt{n}; \lambda_{nj}) \right\} \\ &\leq h_n n^{-1} \sum_{t=1}^n \sum_{i=1}^n \left\{ l(\boldsymbol{\psi}^0(U_t) + \gamma_n \mathbf{w}_t, \mathbf{X}_i, Y_i) - l(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) \right\} K_{h_n}(U_i - U_t) \\ &\quad - h_n n^{-1} \sum_{j=1}^C \sum_{l=1}^{d_j^0} \left\{ p_n(\|\mathbf{b}_{jl}^0 + \gamma_n \mathbf{v}_{jl}\|_2 / \sqrt{n}; \lambda_{nj}) - p_n(\|\mathbf{b}_{jl}^0\|_2 / \sqrt{n}; \lambda_{nj}) \right\} \\ &= D_{n,I}(\mathbf{W}_n) - D_{n,II}(\mathbf{W}_n), \end{aligned} \quad (\text{A.2})$$

where $D_{n,I}(\mathbf{W}_n)$ and $D_{n,II}(\mathbf{W}_n)$ are respectively the differences in the local-kernel log-likelihood and the penalty functions. In what follows, we first perform an order assessment of the two differences for large n , in two steps:

Step 1: (Order assessment of the local log-likelihood difference)

Following the first part of (A.2), we rewrite the difference as

$$D_{n,I}(\mathbf{W}_n) = n^{-1} \sum_{t=1}^n d_{n,I}(U_t, \mathbf{w}_t), \quad (\text{A.3})$$

where for each $t = 1, \dots, n$,

$$d_{n,I}(U_t, \mathbf{w}_t) = h_n \sum_{i=1}^n \{l(\boldsymbol{\psi}^0(U_t) + \gamma_n \mathbf{w}_t, \mathbf{X}_i, Y_i) - l(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i)\} K_{h_n}(U_i - U_t).$$

By the differentiability Condition (RC.2) on f_C^* and using a third-order Taylor expansion,

$$\begin{aligned} d_{n,I}(U_t, \mathbf{w}_t) = & h_n \sum_{i=1}^n \left\{ \gamma_n \mathbf{w}_t^\top l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) + \frac{1}{2} \gamma_n^2 [\mathbf{w}_t^\top l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) \mathbf{w}_t] \right. \\ & \left. + \frac{1}{6} \gamma_n^3 \sum_{j,k,l=1}^{c(d+2)-1} \frac{\partial^3}{\partial \psi_j \partial \psi_k \partial \psi_l} l(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_i, Y_i) w_{j,t} w_{k,t} w_{l,t} \right\} K_{h_n}(U_i - U_t), \end{aligned}$$

where $\tilde{\boldsymbol{\psi}}^0(U_t)$ is between $\boldsymbol{\psi}^0(U_t)$ and $\boldsymbol{\psi}^0(U_t) + \gamma_n \mathbf{w}_t$, for $t = 1, \dots, n$. We then have

$$d_{n,I}(U_t, \mathbf{w}_t) = \mathbf{w}_t^\top l'_{n,h_n}(U_t) + \frac{1}{2} \mathbf{w}_t^\top \{l''_{n,h_n}(U_t)\} \mathbf{w}_t + \frac{h_n \gamma_n^3}{6} \sum_{i=1}^n R(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_i, Y_i), \quad (\text{A.4})$$

where

$$l'_{n,h_n}(U_t) = \sqrt{\frac{h_n}{n}} (1 + r_{1n}) \sum_{i=1}^n l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t), \quad (\text{A.5})$$

$$l''_{n,h_n}(U_t) = \frac{1}{n} (1 + r_{1n})^2 \sum_{i=1}^n l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t)$$

and

$$R(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_i, Y_i) = \sum_{j,k,l=1}^{c(d+2)-1} \frac{\partial^3}{\partial \psi_j \partial \psi_k \partial \psi_l} l(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t) w_{j,t} w_{k,t} w_{l,t}.$$

We now perform an order assessment of the three terms in (A.4), by first focusing on $l'_{n,h_n}(U_t)$. Under our regularity conditions, we will show that $l'_{n,h_n}(U_t) = O_p(1)$. Using Prohorov's Theorem and Example 2.6 on pages 8-10 of van der Vaart (1998), it is sufficient to show that $E_0 \|l'_{n,h_n}(U_t)\|^2 = O(1)$, where E_0 is the expectation with respect to the true joint

distribution of $(U_t, U_i, \mathbf{X}_i, Y_i)$. Note that,

$$\begin{aligned}
& E_0 \|l'_{n,h_n}(U_t)\|^2 \\
&= E_0 \left\{ \left\| (1+r_{1n}) \sqrt{\frac{h_n}{n}} \sum_{i=1}^n l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t) \right\|^2 \right\} \\
&= (1+r_{1n})^2 \frac{h_n}{n} E_0 \left\{ \left\| \sum_{i=1}^n l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t) \right\|^2 \right\} \\
&= (1+r_{1n})^2 \frac{h_n}{n} \sum_{i=1}^n E_0 \left\{ \|l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i)\|^2 K_{h_n}^2(U_i - U_t) \right\} + 2(1+r_{1n})^2 \frac{h_n}{n} \\
&\quad \times \sum_{i < j} E_0 \left\{ [l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i)]^\top [l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_j, Y_j)] K_{h_n}(U_i - U_t) K_{h_n}(U_j - U_t) \right\} \\
&= (1+r_{1n})^2 \frac{h_n}{n} \sum_{i=1}^n E_0[\mathcal{Q}_{i,t}] + 2(1+r_{1n})^2 \frac{h_n}{n} \sum_{i < j} E_0[\mathcal{P}_{i,j,t}], \tag{A.6}
\end{aligned}$$

where $\mathcal{Q}_{i,t}$ and $\mathcal{P}_{i,j,t}$ are the quadratic and cross-product terms, respectively. The s^{th} element of the vector $l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i)$ is denoted by $l'^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i)$, where $s = 1, 2, \dots, d^* = \mathcal{C}(d+2) - 1$.

We now focus on $E_0[\mathcal{Q}_{i,t}]$ for the cases $i = t$ and $i \neq t$. When $i = t$, by the definition of the conditional information matrix $\mathbf{I}(u_t)$ in (12) of the paper and denoting $I^{(s,s)}(u_t)$ as its $(s, s)^{\text{th}}$ element, we obtain

$$\begin{aligned}
E_0[\mathcal{Q}_{t,t}] &= E_0 \left\{ \|l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_t, Y_t)\|^2 K_{h_n}^2(0) \right\} \\
&= K_{h_n}^2(0) \sum_{s=1}^{d^*} E_0 \left\{ l'^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_t, Y_t) \right\}^2 \\
&= K_{h_n}^2(0) \sum_{s=1}^{d^*} E_U \left\{ E_{(\mathbf{X}, Y) | U_t} \left\{ [l'^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_t, Y_t)]^2 \mid U_t \right\} \right\} \\
&= K_{h_n}^2(0) \sum_{s=1}^{d^*} E_U \left\{ I^{(s,s)}(U_t) \right\}.
\end{aligned}$$

By (RC.4), the expected value in the last sum is finite. Thus, for some constant $M > 0$,

$$E_0[\mathcal{Q}_{t,t}] \leq K_{h_n}^2(0) M d^* = \frac{1}{h_n^2} K^2(0) M d^* = \frac{C_0}{h_n^2}. \tag{A.7}$$

For $i \neq t$, we have

$$\begin{aligned} E_0[\mathcal{Q}_{i,t}] &= E_0 \left\{ \left\| l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) \right\|^2 K_{h_n}^2(U_i - U_t) \right\} \\ &= \sum_{s=1}^{d^*} E_0 \left\{ l^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t) \right\}^2. \end{aligned}$$

Next, we evaluate the expected value inside the sum in the above equation. Note that (\mathbf{X}_i, Y_i, U_i) and U_t are independent when $i \neq t$. Thus,

$$\begin{aligned} &E_0 \left\{ l^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t) \right\}^2 \\ &\leq 2E_0 \left\{ l^{(s)}(\boldsymbol{\psi}^0(U_i), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t) \right\}^2 \\ &\quad + 2E_0 \left\{ [l^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) - l^{(s)}(\boldsymbol{\psi}^0(U_i), \mathbf{X}_i, Y_i)] K_{h_n}(U_i - U_t) \right\}^2 \\ &= \mathcal{T}_{n,1} + \mathcal{T}_{n,2}, \end{aligned}$$

where $\mathcal{T}_{n,1}$ and $\mathcal{T}_{n,2}$ respectively represent the two expectations on the right hand side of the above inequality. Order assessment of each term is given below.

By Condition (RC.4) on $\mathbf{I}(u)$, it follows that for some constant $M > 0$,

$$\begin{aligned} \mathcal{T}_{n,1} &= 2 \int_{\mathcal{U}} \int_{\mathcal{U}} \int_{\mathcal{X}} \int_{\mathcal{Y}} \left\{ l^{(s)}(\boldsymbol{\psi}^0(u_i), \mathbf{x}_i, y_i) \right\}^2 f_C^*(y_i | \boldsymbol{\psi}^0(u_i), \mathbf{x}_i) g(\mathbf{x}_i | u_i) \\ &\quad \times K_{h_n}^2(u_i - u_t) m(u_i) m(u_t) dy_i d\mathbf{x}_i du_i du_t \\ &= 2 \int_{\mathcal{U}} \int_{\mathcal{U}} I^{(s,s)}(u_i) K_{h_n}^2(u_i - u_t) m(u_i) m(u_t) du_i du_t \\ &\leq 2M \int_{\mathcal{U}} \int_{\mathcal{U}} K_{h_n}^2(u_i - u_t) m(u_i) m(u_t) du_i du_t. \end{aligned}$$

Applying the change of variable $u_i^* = (u_i - u_t)/h_n$ to the above integration, we have

$$\mathcal{T}_{n,1} \leq 2M h_n^{-1} \int_{\mathcal{U}} \int_{\mathcal{U}^*} K^2(u_i^*) m(u_t + h_n u_i^*) m(u_t) du_i^* du_t.$$

Using a second-order Taylor expansion,

$$m(u_t + h_n u_i^*) = m(u_t) + h_n u_i^* m'(u_t) + (h_n u_i^*)^2 m''(\tilde{u}_{i,t}), \quad (\text{A.8})$$

where $\tilde{u}_{i,t}$ is between u_t and $u_t + h_n u_i^*$. Using the boundedness Condition (RC.2) of $m(u)$ and its first and second derivatives over \mathcal{U} , and the finite-moment Condition (KC.1) of $K^2(u)$,

for some positive constants C_1 and C_2 , we arrive at

$$\mathcal{T}_{n,1} \leq C_1 h_n^{-1} + C_2 h_n. \quad (\text{A.9})$$

Using a first-order Taylor expansion and by Conditions (RC.2) and (RC.3) for some constants C_3 and C_4 ,

$$\begin{aligned} \mathcal{T}_{n,2} &= 2E_0 \left\{ \left[l^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) - l^{(s)}(\boldsymbol{\psi}^0(U_i), \mathbf{X}_i, Y_i) \right] K_{h_n}(U_i - U_t) \right\}^2 \\ &\leq C_3 E_0 \left\{ M_1(\mathbf{X}_i, Y_i)(U_t - U_i) K_{h_n}(U_i - U_t) \right\}^2 \\ &= C_3 \int_{\mathcal{U}} \int_{\mathcal{U}} \left[\int_{\mathcal{X}} \int_{\mathcal{Y}} M_1^2(\mathbf{x}_i, y_i) f_C^*(y_i | \boldsymbol{\psi}^0(u_i), \mathbf{x}_i) g(\mathbf{x}_i | u_i) dy_i d\mathbf{x}_i \right] (u_i - u_t)^2 \\ &\quad \times K_{h_n}^2(u_i - u_t) m(u_i) m(u_t) du_i du_t \\ &\leq C_3 C_4 \int_{\mathcal{U}} \int_{\mathcal{U}} (u_i - u_t)^2 K_{h_n}^2(u_i - u_t) m(u_i) m(u_t) du_i du_t. \end{aligned}$$

Applying the change of variable $u_i^* = (u_i - u_t)/h_n$ and the expansion (A.8), we have

$$\begin{aligned} \mathcal{T}_{n,2} &\leq C_3 C_4 \left\{ h_n \int_{\mathcal{U}} \int_{\mathcal{U}^*} (u_i^*)^2 K^2(u_i^*) m^2(u_t) du_i^* du_t \right. \\ &\quad + h_n^2 \int_{\mathcal{U}} \int_{\mathcal{U}^*} (u_i^*)^3 K^2(u_i^*) m'(u_t) m(u_t) du_i^* du_t \\ &\quad \left. + h_n^3 \int_{\mathcal{U}} \int_{\mathcal{U}^*} (u_i^*)^4 K^2(u_i^*) m''(\tilde{u}_{i,t}) m(u_t) du_i^* du_t \right\}. \end{aligned}$$

Thus, by Condition (RC.2) on $m(u)$ and its first and second derivatives, and the finite-moment Condition (KC.1) of $K^2(u)$, for some positive constants C_5 and C_6 , we arrive at

$$\mathcal{T}_{n,2} \leq C_5 h_n + C_6 h_n^3. \quad (\text{A.10})$$

Hence, using (A.7), (A.9), and (A.10), the first sum on the right hand side of (A.6) becomes

$$\begin{aligned} (1 + r_{1n})^2 \frac{h_n}{n} \sum_{i=1}^n E_0[\mathcal{Q}_{i,t}] &= (1 + r_{1n})^2 \frac{h_n}{n} \left\{ \frac{C_0}{h_n^2} + (n-1)(C_1 h_n^{-1} + C_2 h_n + C_5 h_n + C_6 h_n^3) \right\} \\ &= (1 + r_{1n})^2 \left\{ \frac{C_0}{n h_n} + \frac{(n-1)}{n} [C_1 + (C_2 + C_5) h_n^2 + C_6 h_n^4] \right\}. \end{aligned}$$

Under Condition (PC.2), we have $h_n \rightarrow 0$ and $n h_n \rightarrow \infty$, as $n \rightarrow \infty$. Thus,

$$(1 + r_{1n})^2 \frac{h_n}{n} \sum_{i=1}^n E_0[\mathcal{Q}_{i,t}] = O\{(1 + r_{1n})^2\}. \quad (\text{A.11})$$

We now assess the second sum in (A.6). The expectation of $\mathcal{P}_{i,j,t}$, for $i > j$, $i \neq t$ and $j \neq t$

is given by

$$\begin{aligned}
E_0[\mathcal{P}_{i,j,t}] &= E_0 \left\{ [l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i)]^\top l'(\boldsymbol{\psi}^0(U_t), \mathbf{X}_j, Y_j) K_{h_n}(U_i - U_t) K_{h_n}(U_j - U_t) \right\} \\
&= \sum_{s=1}^{d^*} E_0 \left\{ l^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) l^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_j, Y_j) K_{h_n}(U_i - U_t) K_{h_n}(U_j - U_t) \right\} \\
&= \sum_{s=1}^{d^*} E_{U_t} \left\{ E_{(Y, \mathbf{X}, U_j, U_i) | U_t} \left\{ l^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) l^{(s)}(\boldsymbol{\psi}^0(U_t), \mathbf{X}_j, Y_j) \right. \right. \\
&\quad \left. \left. \times K_{h_n}(U_i - U_t) K_{h_n}(U_j - U_t) \right\} \mid U_t \right\}. \tag{A.12}
\end{aligned}$$

Conditioning on U_t , since $i \neq j$, (U_i, \mathbf{X}_i, Y_i) and (U_j, \mathbf{X}_j, Y_j) are independent and identically distributed. In what follows, we first evaluate the conditional expectations.

$$\begin{aligned}
&E_{(Y, \mathbf{X}, U_i, U_j) | U_t} \left\{ l^{(s)}(\boldsymbol{\psi}^0(u_t), \mathbf{X}_i, Y_i) l^{(s)}(\boldsymbol{\psi}^0(u_t), \mathbf{X}_j, Y_j) K_{h_n}(U_i - u_t) K_{h_n}(U_j - u_t) \mid U_t = u_t \right\} \\
&= E_{(Y, \mathbf{X}, U_i) | U_t}^2 \left\{ l^{(s)}(\boldsymbol{\psi}^0(u_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u_t) \mid U_t = u_t \right\}, \quad i \neq t.
\end{aligned}$$

Thus,

$$\begin{aligned}
&E_{(Y, \mathbf{X}, U_i) | U_t} \left\{ l^{(s)}(\boldsymbol{\psi}^0(u_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u_t) \mid U_t = u_t \right\} \\
&= \int_{\mathcal{Y}} \int_{\mathcal{X}} \int_{\mathcal{U}} l'(\boldsymbol{\psi}^0(u_t), \mathbf{x}_i, y_i) K_{h_n}(u_i - u_t) f_C^*(y_i | \boldsymbol{\psi}^0(u_t), \mathbf{x}_i) g(\mathbf{x}_i | u_i) m(u_i) du_i d\mathbf{x}_i dy_i \\
&= \int_{\mathcal{U}} \left[\int_{\mathcal{Y}} \int_{\mathcal{X}} l'(\boldsymbol{\psi}^0(u_t), \mathbf{x}_i, y_i) f_C^*(y_i | \boldsymbol{\psi}^0(u_t), \mathbf{x}_i) g(\mathbf{x}_i | u_i) d\mathbf{x}_i dy_i \right] K_{h_n}(u_i - u_t) m(u_i) du_i \\
&= \int_{\mathcal{U}} \Delta(u_i; u_t) K_{h_n}(u_i - u_t) m(u_i) du_i \\
&= \int_{\mathcal{U}} K_{h_n}(u_i - u_t) \Delta(u_i; u_t) m(u_i) du_i,
\end{aligned}$$

where $\Delta(u_i; u_t)$ is given in (13) of the paper. Note that by the first Bartlett identity verified in Lemma 1, $\Delta(u_t; u_t) = 0$. Denote $H(u_i) = \Delta(u_i; u_t) m(u_i)$, then

$$E_{(Y, \mathbf{X}, U_i) | U_t} \left\{ l^{(s)}(\boldsymbol{\psi}^0(u_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u_t) \mid U_t = u_t \right\} = \int_{\mathcal{U}} K_{h_n}(u_i - u_t) H(u_i) du_i. \tag{A.13}$$

By the differentiability conditions of (f_C^*, g, m) in (RC.2), and Condition (RC.4) on Δ , the first and second derivatives of $H(u_i)$ with respect to u_i , evaluated at $u_i = u_t$, are

$$H'(u_t) = \Delta'(u_t; u_t) m(u_t)$$

and

$$H''(u_t) = \Delta''(u_t; u_t)m(u_t) + 2\Delta'(u_t; u_t)m'(u_t),$$

where $\Delta'(\cdot; u_t)$ and $\Delta''(\cdot; u_t)$ are the first and second derivatives of $\Delta(u_i; u_t)$ with respect to u_i . Similarly, the third derivative of $H(u_i)$ with respect to u_i is calculated. Using a third-order Taylor expansion,

$$\begin{aligned} H(u_i) &= \Delta'(u_t; u_t)m(u_t)(u_i - u_t) + \frac{1}{2} \{ \Delta''(u_t; u_t)m(u_t) + 2\Delta'(u_t; u_t)m'(u_t) \} (u_i - u_t)^2 \\ &\quad + \frac{1}{6} \{ \Delta'''(\tilde{u}_t; u_t)m(\tilde{u}_t) + \Delta''(\tilde{u}_t; u_t)m'(\tilde{u}_t) + 2\Delta'(\tilde{u}_t; u_t)m''(\tilde{u}_t) + 2\Delta'(\tilde{u}_t; u_t)m''(\tilde{u}_t) \} (u_i - u_t)^3, \end{aligned}$$

where \tilde{u}_t is between u_t and u_i , and $\Delta'''(\cdot; u_t)$ is the third derivative of $\Delta(u_i; u_t)$ with respect to u_i . By replacing the above expansion in (A.13), and Condition (KC.2) on the kernel $K(\cdot)$,

$$\begin{aligned} &E_{(Y, \mathbf{X}, U_i) | U_t} \{ l^{(s)}(\boldsymbol{\psi}^0(u_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u_t) | U_t = u_t \} \\ &= \frac{1}{2} \mathcal{K}_2 \{ \Delta''(u_t; u_t)m(u_t) + 2\Delta'(u_t; u_t)m'(u_t) \} h_n^2 \{ 1 + o(1) \}. \end{aligned} \quad (\text{A.14})$$

Thus,

$$\begin{aligned} E_0[\mathcal{P}_{i,j,t}] &= \{ 1 + o(1) \}^2 \frac{h_n^4 \mathcal{K}_2^2}{4} \sum_{s=1}^{d^*} \int_{\mathcal{U}} \{ \Delta''(u_t; u_t)m(u_t) + 2\Delta'(u_t; u_t)m'(u_t) \}^2 m(u_t) du_t \\ &= \{ 1 + o(1) \}^2 \frac{d^* h_n^4 \mathcal{K}_2^2}{4} \int_{\mathcal{U}} \{ \Delta''(u_t; u_t)m(u_t) + 2\Delta'(u_t; u_t)m'(u_t) \}^2 m(u_t) du_t. \end{aligned}$$

Hence, for $i > j$, $i \neq t$ and $j \neq t$, by the boundedness Condition (RC.2) of $m(u)$ and its first derivative over \mathcal{U} , and Condition (RC.4) on Δ , it follows that, as $n \rightarrow \infty$,

$$E_0[\mathcal{P}_{i,j,t}] = O(h_n^4). \quad (\text{A.15})$$

For $i > j$, and $i = t$, we have

$$\begin{aligned} &E_0[\mathcal{P}_{i,j,i}] \\ &= E_0 \{ [l'(\boldsymbol{\psi}^0(U_i), \mathbf{X}_i, Y_i)]^\top l'(\boldsymbol{\psi}^0(U_i), \mathbf{X}_j, Y_j) K_{h_n}(0) K_{h_n}(U_j - U_i) \} \\ &= K_{h_n}(0) \sum_{s=1}^{d^*} E_{U_i} \{ E_{(Y, \mathbf{X}, U_j) | U_i} \{ l^{(s)}(\boldsymbol{\psi}^0(U_i), \mathbf{X}_i, Y_i) l^{(s)}(\boldsymbol{\psi}^0(U_i), \mathbf{X}_j, Y_j) K_{h_n}(U_j - U_i) | U_i \} \}. \end{aligned}$$

Now, we calculate the conditional expectations in the above sum. By the independence of

(\mathbf{X}_i, Y_i) and (U_j, \mathbf{X}_j, Y_j) for $i \neq j$, and the first Bartlett identity verified in Lemma 1,

$$\begin{aligned}
& E_{(Y, \mathbf{X}, U_j)|U_i} \{l^{(s)}(\boldsymbol{\psi}^0(u_i), \mathbf{X}_i, Y_i)l^{(s)}(\boldsymbol{\psi}^0(u_i), \mathbf{X}_j, Y_j)K_{h_n}(U_j - u_i)|U_i = u_i\} \\
&= E_{(Y, \mathbf{X})|U} \{l^{(s)}(\boldsymbol{\psi}^0(u_i), \mathbf{X}_i, Y_i)|U_i = u_i\} \\
&\quad \times E_{(Y, \mathbf{X}, U_j)|U_i} \{l^{(s)}(\boldsymbol{\psi}^0(u_i), \mathbf{X}_j, Y_j)K_{h_n}(U_j - u_i)|U_i = u_i\} \\
&= 0 \times E_{(Y, \mathbf{X}, U_j)|U_i} \{l^{(s)}(\boldsymbol{\psi}^0(u_i), \mathbf{X}_j, Y_j)K_{h_n}(U_j - u_i)|U_i = u_i\} = 0.
\end{aligned}$$

Thus,

$$E_0[\mathcal{P}_{i,j,i}] = E_0[\mathcal{P}_{i,j,j}] = 0. \quad (\text{A.16})$$

Using (A.15) and (A.16), the second sum in (A.6) can be written as

$$\begin{aligned}
2(1 + r_{1n})^2 \frac{h_n}{n} \sum_{i < j} E_0[\mathcal{P}_{i,j,t}] &= 2(1 + r_{1n})^2 \frac{h_n}{n} \left\{ \sum_{\substack{i < j \\ i, j \neq t}}^n E_0[\mathcal{P}_{i,j,t}] + E_0[\mathcal{P}_{i,j,i}] + E_0[\mathcal{P}_{i,j,j}] \right\} \\
&= 2(1 + r_{1n})^2 \frac{h_n}{n} \left\{ \sum_{\substack{i < j \\ i, j \neq t}}^n O(h_n^4) + 0 + 0 \right\} \\
&= (1 + r_{1n})^2 \frac{(n-1)(n-2)}{n} O(h_n^5) = (1 + r_{1n})^2 O(nh_n^5).
\end{aligned}$$

Under conditions (PC.2), as $n \rightarrow \infty$, we have $nh_n^5 = O(1)$. Thus, for large n

$$2(1 + r_{1n})^2 \frac{h_n}{n} \sum_{i < j} E_0[\mathcal{P}_{i,j,t}] = O\{(1 + r_{1n})^2\}. \quad (\text{A.17})$$

Hence, by putting together the order assessments in (A.11) and (A.17), and using (A.6), for large n we arrive at $E_0\|l'_{n,h_n}(U_t)\|^2 = O\{(1 + r_{1n})^2\}$, which implies that

$$l'_{n,h_n}(U_t) = O_p\{(1 + r_{1n})\}. \quad (\text{A.18})$$

Next, we perform an order assessment of the third term in (A.4). Using Condition (RC.3),

$$\begin{aligned}
& E_0 \left[\left| \frac{h_n \gamma_n^3}{6} \sum_{i=1}^n R \left(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_i, Y_i \right) \right| \right] \\
& \leq \frac{h_n \gamma_n^3}{6} \sum_{i=1}^n E_0 \left| R \left(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_i, Y_i \right) \right| \\
& = \frac{h_n \gamma_n^3}{6} \left[E_0 \left\{ \left| R \left(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_t, Y_t \right) \right| \right\} + \sum_{(t \neq i)=1}^n E_0 \left\{ \left| R \left(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_i, Y_i \right) \right| \right\} \right] \\
& = \frac{h_n \gamma_n^3}{6} K_{h_n}(0) \sum_{j,k,l=1}^{d^*} E_0 \left[\left| \frac{\partial^3}{\partial \psi_j \partial \psi_k \partial \psi_l} l \left(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_t, Y_t \right) w_{j,t} w_{k,t} w_{l,t} \right| \right] \\
& + \frac{h_n \gamma_n^3}{6} E_0 \left[\left| \sum_{(t \neq i)=1}^n \sum_{j,k,l=1}^{d^*} \frac{\partial^3}{\partial \psi_j \partial \psi_k \partial \psi_l} l \left(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_i, Y_i \right) K_{h_n}(U_i - U_t) w_{j,t} w_{k,t} w_{l,t} \right| \right] \\
& \leq \frac{(d^*)^3 h_n \gamma_n^3}{6} K_{h_n}(0) E_0 [M_2(\mathbf{X}_t, Y_t)] \left\{ \sum_{l=1}^{d^*} |w_{l,t}|^2 \right\}^{3/2} \\
& + \frac{(d^*)^3 h_n \gamma_n^3}{6} \left[\sum_{(t \neq i)=1}^n E_0 [M_2(\mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t)] \right] \left\{ \sum_{l=1}^{d^*} |w_{l,t}|^2 \right\}^{3/2} \\
& \leq \frac{(d^*)^3 h_n \gamma_n^3}{6} K_{h_n}(0) E_0 [M_2(\mathbf{X}_t, Y_t)] \left\{ \sum_{l=1}^{d^*} |w_{l,t}|^2 \right\}^{3/2} \\
& + \frac{(d^*)^3 h_n \gamma_n^3}{6} n E_0 [M_2(\mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t)] \left\{ \sum_{l=1}^{d^*} |w_{l,t}|^2 \right\}^{3/2} \\
& = \frac{(d^*)^3 \gamma_n^3}{6} K(0) E_0 [M_2(\mathbf{X}_t, Y_t)] \left\{ \sum_{j=1}^{d^*} |w_{j,t}|^2 \right\}^{3/2} \\
& + \frac{(d^*)^3 h_n \gamma_n^3}{6} n E_0 [M_2(\mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t)] \left\{ \sum_{l=1}^{d^*} |w_{l,t}|^2 \right\}^{3/2} \\
& = O(\gamma_n^3) + O(\gamma_n) = O(\gamma_n).
\end{aligned}$$

Thus,

$$\frac{h_n \gamma_n^3}{6} \sum_{i=1}^n R \left(\tilde{\boldsymbol{\psi}}^0(U_t), \mathbf{X}_i, Y_i \right) = O_p(\gamma_n) \tag{A.19}$$

uniformly for all $t = 1, 2, \dots, n$. Thus, since $\gamma_n \rightarrow 0$, (A.4) reduces to

$$d_{n,I}(U_t, \mathbf{w}_t) = \mathbf{w}_t^\top l'_{n,h_n}(U_t) + \frac{1}{2} \mathbf{w}_t^\top l''_{n,h_n}(U_t) \mathbf{w}_t + o_p(1). \tag{A.20}$$

On the other hand,

$$\begin{aligned}
l''_{n,h_n}(U_t) &= \frac{(1+r_{1n})^2}{n} \sum_{i=1}^n l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t) \\
&= \frac{K(0)}{nh} (1+r_{1n})^2 l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_t, Y_t) \\
&\quad + \frac{(1+r_{1n})^2}{n} \left\{ \sum_{(t \neq) i=1}^n l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_i, Y_i) K_{h_n}(U_i - U_t) \right\} \\
&= \mathcal{E}_{1n} + \mathcal{E}_{2n}(U_t).
\end{aligned} \tag{A.21}$$

Note that

$$-E_0 \{l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_t, Y_t)\} = E_U \{\mathbf{I}(U_t)\}$$

which by Condition (RC.4) is finite and positive definite. Thus, by Condition (RC.5),

$$\begin{aligned}
&l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_t, Y_t) \\
&= \{l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_t, Y_t) - E_0 \{l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_t, Y_t)\}\} + E_0 \{l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_t, Y_t)\} \\
&= O_p(1) - E_0 \{\mathbf{I}(U_t)\} = O_p(1).
\end{aligned}$$

Hence, by conditions $r_{1n} = O(1)$ and $nh_n \rightarrow \infty$ (see (PC.2)), as $n \rightarrow \infty$,

$$\mathcal{E}_{1n} = \frac{K(0)}{nh} (1+r_{1n})^2 l''(\boldsymbol{\psi}^0(U_t), \mathbf{X}_t, Y_t) = \frac{(1+r_{1n})^2}{nh_n} O_p(1) = o_p(1). \tag{A.22}$$

Next, we focus on $\mathcal{E}_{2n}(U_t)$. Define the centralized (finite dimensional) matrix

$$Z_{i,n}(u) = l''(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u) - E_0[l''(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u)], \tag{A.23}$$

where $E_0[Z_{i,n}(u)] = 0$, for any $u \in \mathcal{U}$. Note that, by using a change of variable $\nu_i = (u_i - u)/h_n$,

$$\begin{aligned}
&E_0[l''(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u)] \\
&= \int_{\mathcal{Y}} \int_{\mathcal{X}} \int_{\mathcal{U}} l''(\boldsymbol{\psi}^0(u), \mathbf{x}_i, y_i) K_{h_n}(u_i - u) f(u_i, \mathbf{x}_i, y_i) du_i d\mathbf{x}_i dy_i \\
&= \int_{\mathcal{U}} \int_{\mathcal{Y}} \int_{\mathcal{X}} l''(\boldsymbol{\psi}^0(u), \mathbf{x}_i, y_i) K(\nu_i) f(u + h_n \nu_i, \mathbf{x}_i, y_i) d\mathbf{x}_i dy_i d\nu_i,
\end{aligned}$$

where f is the joint pdf of (U_i, \mathbf{X}_i, Y_i) . By a second-order Taylor expansion,

$$f(u + h_n \nu_i, \mathbf{x}_i, y_i) = f(u, \mathbf{x}_i, y_i) + h_n \nu_i f'(u, \mathbf{x}_i, y_i) + \frac{h_n^2 \nu_i^2}{2} f''(u, \mathbf{x}_i, y_i) \tag{A.24}$$

such that u_{ni} lies between u and $u + h_n\nu_i$, and f' and f'' are the partial derivatives of $f(u, \mathbf{x}_i, y_i)$ with respect to u , which by Conditions (RC.1) and (RC.2) exist. Replacing (A.24) in the above integration and using the boundedness Condition (KC.2) on the kernel $K(\cdot)$ and by the second Bartlett identity verified in Lemma 1, we arrive at

$$E_0[l''(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i)K_{h_n}(U_i - u)] = -m(u)\mathbf{I}(u) + o_p(h_n), \quad \forall u \in \mathcal{U}. \quad (\text{A.25})$$

Turning to (A.23), we show that, as $n \rightarrow \infty$,

$$\sup_{u \in \mathcal{U}} \left\| \frac{1}{n} \sum_{i=1}^n Z_{i,n}(u) \right\|_2 = O_p \left\{ \left(\frac{\log(\frac{1}{h_n})}{nh_n} \right)^{\frac{1}{2}} \right\} \equiv O_p(\xi_n) = o_p(1). \quad (\text{A.26})$$

Let $Z_{i,n}^{jl}(u)$ be the (j, l) -th element of the matrix $Z_{i,n}(u)$. Since \mathcal{U} is compact, let $\bigcup_{k=1}^{\mathcal{N}_n} B(u_k; \eta_n)$ be a finite open cover of \mathcal{U} , where $B(u_k; \eta_n)$ are balls centred at u_k and with radius η_n , such that $\eta_n = o(1)$, as $n \rightarrow \infty$. Clearly, we have $\mathcal{N}_n = O(\eta_n^{-1})$. Then, we can write,

$$\begin{aligned} \sup_{u \in \mathcal{U}} \left| \frac{1}{n} \sum_{i=1}^n Z_{i,n}^{jl}(u) \right| &\leq \max_{1 \leq k \leq \mathcal{N}_n} \left| \frac{1}{n} \sum_{i=1}^n Z_{i,n}^{jl}(u_k) \right| + \max_{1 \leq k \leq \mathcal{N}_n} \sup_{u \in \mathcal{U}} \left| \frac{1}{n} \sum_{i=1}^n [Z_{i,n}^{jl}(u) - Z_{i,n}^{jl}(u_k)] \right| \\ &\leq \max_{1 \leq k \leq \mathcal{N}_n} \left| \frac{1}{n} \sum_{i=1}^n Z_{i,n}^{jl}(u_k) \right| + \max_{1 \leq k \leq \mathcal{N}_n} \sup_{u \in B(u_k; \eta_n)} \left| \frac{1}{n} \sum_{i=1}^n [Z_{i,n}^{jl}(u) - Z_{i,n}^{jl}(u_k)] \right| \\ &= \mathcal{Z}_{n1} + \mathcal{Z}_{n2}. \end{aligned} \quad (\text{A.27})$$

We show that \mathcal{Z}_{n1} and \mathcal{Z}_{n2} follow the rate as in (A.26). We first focus on \mathcal{Z}_{n2} . Note that,

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n [Z_{i,n}^{jl}(u) - Z_{i,n}^{jl}(u_k)] \\ &= \frac{1}{n} \sum_{i=1}^n \{ [l''_{jl}(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i)K_{h_n}(U_i - u) - l''_{jl}(\boldsymbol{\psi}^0(u_k), \mathbf{X}_i, Y_i)K_{h_n}(U_i - u_k)] \\ &\quad - E_0 [l''_{jl}(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i)K_{h_n}(U_i - u) - l''_{jl}(\boldsymbol{\psi}^0(u_k), \mathbf{X}_i, Y_i)K_{h_n}(U_i - u_k)] \}, \end{aligned}$$

where l''_{jl} is the (j, l) -th element of the matrix l'' . The expression in the first square bracket can be rewritten as

$$\begin{aligned} &[l''_{jl}(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i)K_{h_n}(U_i - u) - l''_{jl}(\boldsymbol{\psi}^0(u_k), \mathbf{X}_i, Y_i)K_{h_n}(U_i - u_k)] \\ &= l''_{jl}(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) [K_{h_n}(U_i - u) - K_{h_n}(U_i - u_k)] \\ &\quad + [l''_{jl}(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) - l''_{jl}(\boldsymbol{\psi}^0(u_k), \mathbf{X}_i, Y_i)] K_{h_n}(U_i - u_k). \end{aligned}$$

By the Lipschitz condition in (KC.1) on the kernel K , and Conditions (RC.3) and (RC.5) on $l(\boldsymbol{\psi}(u), \boldsymbol{x}, y)$, we have that, for large n ,

$$\begin{aligned} & \left| \frac{1}{n} \sum_{i=1}^n l''_{jl}(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) [K_{h_n}(U_i - u) - K_{h_n}(U_i - u_k)] \right| \\ & \leq \left| \frac{1}{n} \sum_{i=1}^n l''_{jl}(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) \right| \frac{C|u - u_k|}{h_n^2} = O_p\left(\frac{\eta_n}{h_n^2}\right). \end{aligned}$$

Similarly, by the above conditions and also Condition (RC.2), the Condition (KC.1) on the kernel $K(\cdot)$, and using the mean value theorem, for large n ,

$$\frac{1}{n} \sum_{i=1}^n \{ |l''_{jl}(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) - l''_{jl}(\boldsymbol{\psi}^0(u_k), \mathbf{X}_i, Y_i)| K_{h_n}(U_i - u_k) \} \leq O_p(1) \frac{|u - u_k|}{h_n} = O_p\left(\frac{\eta_n}{h_n}\right).$$

Hence, by the putting together the above order assessments, we have that, for large n ,

$$\left| \frac{1}{n} \sum_{i=1}^n [Z_{i,n}^{jl}(u) - Z_{i,n}^{jl}(u_k)] \right| = O_p\left(\frac{r_n}{h_n^2}\right) + O_p\left(\frac{\eta_n}{h_n}\right) = O_p\left(\frac{\eta_n}{h_n^2}\right)$$

which implies that if we choose $\eta_n = h_n^2 \xi_n$, we have, for large n ,

$$\mathcal{Z}_{n2} = O_p\left(\frac{\eta_n}{h_n^2}\right) = O_p(\xi_n). \quad (\text{A.28})$$

Next, we focus on the order assessment of \mathcal{Z}_{n1} in (A.27). Denote the random variables,

$$\bar{Z}_{i,n}^{jl}(u) = Z_{i,n}^{jl}(u) \mathbf{1}\{|l''_{jl}(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i)| \leq C_2 n^{1/m_0}\} \quad (\text{A.29})$$

and

$$\tilde{Z}_{i,n}^{jl}(u) = Z_{i,n}^{jl}(u) \mathbf{1}\{|l''_{jl}(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i)| > C_2 n^{1/m_0}\} \quad (\text{A.30})$$

such that $Z_{i,n}^{jl}(u) = \bar{Z}_{i,n}^{jl}(u) + \tilde{Z}_{i,n}^{jl}(u)$, for some constant $C_2 > 0$, and $\mathbf{1}\{\cdot\}$ is indicator function.

Note that $E_0\{Z_{i,n}^{jl}(u)\} = 0$. Then, we have that

$$\begin{aligned} \mathcal{Z}_{n1} &= \max_{1 \leq k \leq \mathcal{N}_n} \left| \frac{1}{n} \sum_{i=1}^n Z_{i,n}^{jl}(u_k) \right| \\ &\leq \max_{1 \leq k \leq \mathcal{N}_n} \left| \frac{1}{n} \sum_{i=1}^n [\bar{Z}_{i,n}^{jl}(u_k) - E\{\bar{Z}_{i,n}^{jl}(u_k)\}] \right| + \max_{1 \leq k \leq \mathcal{N}_n} \left| \frac{1}{n} \sum_{i=1}^n [\tilde{Z}_{i,n}^{jl}(u_k) - E\{\tilde{Z}_{i,n}^{jl}(u_k)\}] \right| \\ &= \mathcal{J}_{n1} + \mathcal{J}_{n2}. \end{aligned} \quad (\text{A.31})$$

For some constant $C_3 > 0$, any $\epsilon > 0$, by using the Markov inequality and Condition (RC.3),

$$\begin{aligned} P(\mathcal{J}_{n2} > C_3 \xi_n) &\leq P\left(\max_{1 \leq k \leq \mathcal{N}_n} \bigcup_{i=1}^n \{|l''_{jl}(\boldsymbol{\psi}^0(u_k), \mathbf{X}_i, Y_i)| > C_2 n^{1/m_0}\}\right) \\ &\leq P\left(\bigcup_{i=1}^n \{|M_1(\mathbf{X}_i, Y_i)| > C_2 n^{1/m_0}\}\right) \leq \sum_{i=1}^n P(M_1(\mathbf{X}_i, Y_i) > C_2 n^{1/m_0}) \\ &\leq \sum_{i=1}^n E_0\{M_1^{m_0}(\mathbf{X}_i, Y_i)\} (C_2 n^{1/m_0})^{-m_0} = E_0\{M_1^{m_0}(\mathbf{X}_i, Y_i)\} C_2^{-m_0} < \epsilon \end{aligned}$$

if we choose $C_2 > \epsilon^{-1/m_0} [E_0\{M_1^{m_0}(\mathbf{X}_i, Y_i)\}]^{1/m_0}$. Since $\epsilon > 0$ can be arbitrarily small, we then have that, for large n ,

$$\mathcal{J}_{n2} = O_p(\xi_n). \quad (\text{A.32})$$

We use the Bernstein inequality, Lemma 2.2.9 of [van der Vaart and Wellner \(1996\)](#), to assess the large sample behavior of \mathcal{J}_{n1} in (A.31). First, we verify the conditions of the lemma. Note that, by the definitions in (A.23) and (A.29), and the boundedness of $K(\cdot)$,

$$\left| \bar{Z}_{i,n}^{jl}(u_k) - E\{\bar{Z}_{i,n}^{jl}(u_k)\} \right| \leq \frac{C_4 C_2 n^{1/m_0}}{h_n}, \quad \text{Var}\{\bar{Z}_{i,n}^{jl}(u_k)\} \leq \frac{C_4}{h_n} \quad (\text{A.33})$$

for some constant $C_4 > 0$. By the Bernstein inequality,

$$P(\mathcal{J}_{n1} > C_3 \xi_n) \leq 2\mathcal{N}_n \exp\left\{-\frac{C_3^2 n^2 \xi_n^2}{2nC_4/h_n + 2nC_4 \xi_n C_3 C_2 n^{1/m_0}/3h_n}\right\}.$$

If we choose C_3 such that $C_3 > 3C_4$, then by the definition of ξ_n ,

$$\frac{C_3^2 n^2 \xi_n^2}{2nC_4/h_n + 2nC_4 \xi_n C_3 C_2 n^{1/m_0}/3h_n} \geq \frac{(3C_3/2) \log h_n^{-1}}{1 + C_2 C_4 \xi_n n^{1/m_0}}.$$

Also if, for large n ,

$$n^{1/m_0} \xi_n = \left(\frac{n^{2/m_0} \log(\frac{1}{h_n})}{nh_n}\right)^{\frac{1}{2}} = o(1), \quad \mathcal{N}_n \exp\{-C_3 \log h_n^{-1}\} = o(1) \quad (\text{A.34})$$

Thus, for large n , we have that $P(\mathcal{J}_{n1} > C_3 \xi_n) \leq o(1)$, which implies that

$$\mathcal{J}_{n1} = O_p(\xi_n). \quad (\text{A.35})$$

Hence, in summary, by (A.28), and also applying (A.32) and (A.35) to (A.31), we have proved (A.26).

Therefore, for any $u \in \mathcal{U}$, and using (A.25),

$$\begin{aligned} & \frac{1}{n} \left\{ \sum_{i=1}^n l''(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u) \right\} \\ &= \frac{1}{n} \sum_{i=1}^n Z_{i,n}(u) + E_0 [l''(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u)] \\ &= \frac{1}{n} \sum_{i=1}^n Z_{i,n}(u) - m(u) \mathbf{I}(u) + o_p(h_n) \end{aligned} \quad (\text{A.36})$$

and then, using (A.26), we have that

$$\begin{aligned} & \sup_{u \in \mathcal{U}} \left\| \frac{1}{n} \sum_{i=1}^n l''(\boldsymbol{\psi}^0(u), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u) + m(u) \mathbf{I}(u) \right\|_2 \\ &= \sup_{u \in \mathcal{U}} \left\| \frac{1}{n} \sum_{i=1}^n Z_{i,n}(u) \right\|_2 + o_p(h_n) = O_p(\xi_n + h_n). \end{aligned} \quad (\text{A.37})$$

This implies that, for large n ,

$$\begin{aligned} & \sup_{u \in \mathcal{U}} \|l''_{n,h_n}(u) + (1 + r_{1n})^2 m(u) \mathbf{I}(u)\|_2 \\ &= \sup_{u \in \mathcal{U}} \|o_p\{(1 + r_{1n})^2\} + \mathcal{E}_{2n}(u) + (1 + r_{1n})^2 m(u) \mathbf{I}(u)\|_2 \\ &= (1 + r_{1n})^2 O_p(\xi_n + h_n) = o_p\{(1 + r_{1n})^2\} \end{aligned} \quad (\text{A.38})$$

since, by Condition (PC.2), $nh_n \rightarrow \infty$.

Using (A.20) and taking into account (A.38), then (A.3) reduces to

$$\begin{aligned} D_{n,I}(\mathbf{W}_n) &= \frac{1}{n} \sum_{t=1}^n \left[\mathbf{w}_t^\top l'_{n,h_n}(U_t) + \frac{1}{2} \mathbf{w}_t^\top l''_{n,h_n}(U_t) \mathbf{w}_t \right] + o_p(1) \\ &= \frac{1}{n} \sum_{t=1}^n \left[\mathbf{w}_t^\top l'_{n,h_n}(U_t) + \frac{1}{2} \mathbf{w}_t^\top \{l''_{n,h_n}(U_t) + (1 + r_{1n})^2 m(U_t) \mathbf{I}(U_t) - (1 + r_{1n})^2 m(U_t) \mathbf{I}(U_t)\} \mathbf{w}_t \right] \\ &\quad + o_p(1) \\ &\leq \frac{1}{n} \sum_{t=1}^n \left[\mathbf{w}_t^\top l'_{n,h_n}(U_t) - \frac{(1 + r_{1n})^2}{2} \mathbf{w}_t^\top \{m(U_t) \mathbf{I}(U_t) + o_p(1)\} \mathbf{w}_t \right] + o_p(1) \\ &= \frac{1}{n} \sum_{t=1}^n \left[\mathbf{w}_t^\top l'_{n,h_n}(U_t) - \frac{(1 + r_{1n})^2}{2} \mathbf{w}_t^\top \{m(U_t) \mathbf{I}(U_t)\} \mathbf{w}_t \right] + o_p(1). \end{aligned} \quad (\text{A.39})$$

By conditions (RC.2) and (RC.4), the matrix $m(U_t) \mathbf{I}(U_t)$ is positive definite, and thus all of its eigenvalues are positive. Let λ_t^{\min} be the smallest eigenvalue of $m(U_t) \mathbf{I}(U_t)$. Also, let $\lambda_{\min} = \min\{\lambda_t^{\min}, t = 1, \dots, n\}$ and $\lambda_0^{\min} = \inf_{u \in \mathcal{U}} \lambda_{\min}\{m(u) \mathbf{I}(u)\}$. Inequality in (A.39)

reduces to

$$\begin{aligned}
& D_{n,I}(\mathbf{W}_n) \\
& \leq \frac{1}{n} \sum_{t=1}^n \left[\mathbf{w}_t^\top l'_{n,h_n}(U_t) - \frac{(1+r_{1n})^2}{2} \mathbf{w}_t^\top \lambda_t^{\min} \mathbf{w}_t \right] + o_p(1) \\
& \leq \frac{1}{n} \left[\sum_{t=1}^n \mathbf{w}_t^\top l'_{n,h_n}(U_t) - \frac{(1+r_{1n})^2}{2} \times \lambda_{\min} \times \sum_{t=1}^n \|\mathbf{w}_t\|^2 \right] + o_p(1) \\
& \leq (n^{-1} \|\mathbf{W}_n\|^2)^{\frac{1}{2}} \left(n^{-1} \sum_{t=1}^n \|l'_{n,h_n}(U_t)\|^2 \right)^{\frac{1}{2}} - \frac{(1+r_{1n})^2}{2} \times \lambda_{\min} \times (n^{-1} \|\mathbf{W}_n\|^2) + o_p(1) \\
& = M_\epsilon \left(n^{-1} \sum_{t=1}^n \|l'_{n,h_n}(U_t)\|^2 \right)^{\frac{1}{2}} - \frac{(1+r_{1n})^2}{2} \times \lambda_{\min} \times M_\epsilon^2 + o_p(1) \\
& \leq M_\epsilon \left(n^{-1} \sum_{t=1}^n \|l'_{n,h_n}(U_t)\|^2 \right)^{\frac{1}{2}} - \frac{(1+r_{1n})^2}{2} \times \lambda_0^{\min} \times M_\epsilon^2 (1 + o_p(1)). \tag{A.40}
\end{aligned}$$

As we showed in (A.18), we have

$$n^{-1} \sum_{t=1}^n E \|l'_{n,h_n}(U_t)\|^2 = n^{-1} \sum_{t=1}^n O_p\{(1+r_{1n})^2\} = O_p\{(1+r_{1n})^2\}$$

Note that $\lambda_0^{\min} > 0$. Therefore, we can choose M_ϵ large enough such that the second term in (A.40) and thus $D_{n,I}(\mathbf{W}_n)$ become negative, in probability, for large n .

Step 2: (Order assessment of the penalty difference) Following the penalty difference

in (A.2), $D_{n,II}(\mathbf{W}_n)$, which involves the grouping penalty and using Taylor expansion,

$$\begin{aligned}
& D_{n,II}(\mathbf{W}_n) \tag{A.41} \\
&= h_n n^{-1} \sum_{j=1}^c \sum_{l=1}^{d_j^0} \left\{ p_n(\|\mathbf{b}_{jl}^0 + \gamma_n \mathbf{v}_{jl}\|/\sqrt{n}; \lambda_{nj}) - p_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj}) \right\} \\
&= h_n n^{-1} \sum_{j=1}^c \sum_{l=1}^{d_j^0} \sum_{t=1}^n \frac{\gamma_n \beta_{jl}^0(U_t)}{\sqrt{n} \|\mathbf{b}_{jl}^0\|} p'_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj}) w_{jl,I,t} \\
&\quad + h_n n^{-1} \sum_{j=1}^c \sum_{l=1}^{d_j^0} \sum_{t=1}^n \frac{\gamma_n^2}{2\sqrt{n}} \left(\frac{1}{\|\mathbf{b}_{jl}^0\|} - \frac{[\beta_{jl}^0(U_t)]^2}{\|\mathbf{b}_{jl}^0\|^3} \right) p'_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj}) w_{cj,I,t}^2 \\
&\quad + h_n n^{-1} \sum_{j=1}^c \sum_{l=1}^{d_j^0} \sum_{t=1}^n \frac{\gamma_n^2 [\beta_{cj}^0(U_t)]^2}{2n \|\mathbf{b}_{jl}^0\|^2} p''_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj}) w_{jl,I,t}^2 \\
&\quad + h_n n^{-1} \sum_{j=1}^c \sum_{l=1}^{d_j^0} \sum_{t \neq l=1}^n \frac{-\gamma_n^2 \beta_{jl}^0(U_t) \beta_{jl}^0(U_l)}{\sqrt{n} \|\mathbf{b}_{jl}^0\|^3} p'_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj}) w_{jl,I,t} w_{jl,I,l}/2 \\
&\quad + h_n n^{-1} \sum_{j=1}^c \sum_{l=1}^{d_j^0} \sum_{t \neq l=1}^n \frac{\gamma_n^2 \beta_{jl}^0(U_t) \beta_{jl}^0(U_l)}{n \|\mathbf{b}_{jl}^0\|^2} p''_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj}) w_{jl,I,t} w_{jl,I,l}/2 \\
&= d_{1,n,II}(\mathbf{W}_n) + d_{2,n,II}(\mathbf{W}_n) + d_{3,n,II}(\mathbf{W}_n) + d_{4,n,II}(\mathbf{W}_n) + d_{5,n,II}(\mathbf{W}_n).
\end{aligned}$$

We now perform order assessment of the terms $d_{1,n,II}(\mathbf{w}), \dots, d_{5,n,II}(\mathbf{w})$. For large n ,

$$\begin{aligned}
& d_{1,n,II}(\mathbf{W}_n) \\
&\leq (1 + r_{1n}) h_n^{\frac{1}{2}} n^{-\frac{4}{2}} \sum_{j=1}^c \sum_{l=1}^{d_j^0} \sum_{t=1}^n \frac{|\beta_{jl}^0(U_t)|}{\|\mathbf{b}_{jl}^0\|} \max_{\substack{1 \leq l \leq d_j^0 \\ 1 \leq j \leq c}} |p'_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj})| |w_{jl,I,t}| \\
&\leq (1 + r_{1n}) h_n^{\frac{1}{2}} \max_{\substack{1 \leq l \leq d_j^0 \\ 1 \leq j \leq c}} \left| \frac{p'_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj})}{n^{3/2}} \right| \sum_{j=1}^c \sum_{l=1}^{d_j^0} \left\{ \sum_{t=1}^n w_{jl,I,t}^2/n \right\}^{1/2} \\
&\leq M_\epsilon (1 + r_{1n}) r_{1n}. \tag{A.42}
\end{aligned}$$

Following $d_{2,n,II}(\mathbf{W}_n)$, for large n ,

$$\begin{aligned}
& |d_{2,n,II}(\mathbf{W}_n)| \\
& \leq 2^{-1} h_n n^{-3/2} \gamma_n^2 \sum_{j=1}^{\mathcal{C}} \sum_{l=1}^{d_\kappa^0} \sum_{t=1}^n \frac{1}{\|\mathbf{b}_{jl}^0\|} \left(1 + \frac{[\beta_{jl}^0(U_t)]^2}{\|\mathbf{b}_{jl}^0\|^2} \right) |p'_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj})| w_{jl,I,t}^2 \\
& \leq (1+r_{1n})^2 n^{-1/2} \left\{ \min_{j,l} \|\mathbf{b}_{jl}^0\|/\sqrt{n} \right\}^{-1} \max_{\substack{1 \leq l \leq d_j^0 \\ 1 \leq j \leq \mathcal{C}}} \left| \frac{p'_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj})}{n^{3/2}} \right| \left\{ \sum_{t=1}^n w_{jl,I,t}^2/n \right\} \\
& = M_\epsilon^2 r_{1n} (1+r_{1n})^2 \{nh_n\}^{-1/2} \left\{ \min_{j,l} E_0 [\beta_{jl}^2(U)] \right\}^{-1/2}. \tag{A.43}
\end{aligned}$$

Following $d_{3,n,II}(\mathbf{W}_n)$, for large n ,

$$\begin{aligned}
|d_{3,n,II}(\mathbf{W}_n)| & \leq h_n n^{-1} \sum_{j=1}^{\mathcal{C}} \sum_{l=1}^{d_j^0} \sum_{t=1}^n \frac{\gamma_n^2 [\beta_{jl}^0(U_t)]^2}{2n \|\mathbf{b}_{jl}^0\|^2} |p''_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj})| w_{jl,I,t}^2 \\
& \leq (1+r_{1n})^2 \max_{\substack{1 \leq l \leq d_j^0 \\ 1 \leq j \leq \mathcal{C}}} \left| \frac{p''_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj})}{n^2} \right| \left\{ \sum_{t=1}^n w_{jl,I,t}^2/n \right\} \\
& \leq M_\epsilon^2 r_{2n} (1+r_{1n})^2. \tag{A.44}
\end{aligned}$$

Following $d_{4,n,II}(\mathbf{W}_n)$,

$$\begin{aligned}
& |d_{4,n,II}(\mathbf{W}_n)| \\
& \leq h_n n^{-1} \sum_{j=1}^{\mathcal{C}} \sum_{l=1}^{d_\kappa^0} \sum_{t \neq l=1}^n \frac{\gamma_n^2 |\beta_{jl}^0(U_t)| |\beta_{jl}^0(U_l)|}{\sqrt{n} \|\mathbf{b}_{jl}^0\|^3} |p'_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{n\kappa})| |w_{jl,I,t} w_{jl,I,l}| / 2 \\
& \leq h_n n^{-3/2} (1+r_{1n})^2 (h_n n)^{-1} \left\{ \min_{j,l} \|\mathbf{b}_{jl}^0\| \right\}^{-1} \max_{\substack{1 \leq j \leq d_\kappa^0 \\ 1 \leq \kappa \leq \mathcal{C}}} |p'_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj})| \left\{ \sum_{t=1}^n w_{jl,I,t}^2 \right\} \\
& = n^{-1/2} (1+r_{1n})^2 \left\{ \min_{j,l} \|\mathbf{b}_{jl}^0\|/\sqrt{n} \right\}^{-1} \max_{\substack{1 \leq l \leq d_j^0 \\ 1 \leq j \leq \mathcal{C}}} \left| \frac{p'_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj})}{n^{3/2}} \right| \left\{ \sum_{t=1}^n w_{jl,I,t}^2/n \right\} \\
& = M_\epsilon^2 \{nh_n\}^{-1/2} r_{1n} (1+r_{1n})^2 \left\{ \min_{j,l} E_0 [\beta_{jl}^2(U)] \right\}^{-1/2}. \tag{A.45}
\end{aligned}$$

Following $d_{5,n,II}(\mathbf{W}_n)$,

$$\begin{aligned}
& |d_{5,n,II}(\mathbf{W}_n)| \\
& \leq h_n n^{-1} \sum_{j=1}^c \sum_{l=1}^{d_j^0} \sum_{t \neq l=1}^n \frac{\gamma_n^2 |\beta_{jl}^0(U_t)| |\beta_{jl}^0(U_l)|}{n \|\mathbf{b}_{jl}^0\|^2} p_n''(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj}) |w_{jl,I,t} w_{jl,I,l}|/2 \\
& \leq (1 + r_{1n})^2 \max_{\substack{1 \leq l \leq d_j^0 \\ 1 \leq j \leq c}} \left| \frac{p_n''(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj})}{n^2} \right| \left\{ \sum_{t=1}^n w_{jl,I,t}^2/n \right\} \\
& = M_\epsilon^2 r_{2n} (1 + r_{1n})^2. \tag{A.46}
\end{aligned}$$

Thus, by Condition (PC.2) on the penalty p_n and the smoothing parameter h_n , the order assessment in (A.40) and those in (A.42)-(A.46) corresponding to the terms $D_{n,I}(\mathbf{W}_n)$ and $D_{n,II}(\mathbf{W}_n)$ in (A.2) imply that for a sufficiently large M_ϵ the expression

$$-\frac{(1 + r_{1n})^2}{2} \lambda_0^{\min} M_\epsilon^2$$

is the sole leading term in the right side of $D_n(\mathbf{w})$ in (A.2). Therefore, for any given $\epsilon > 0$, there exists a sufficiently large M_ϵ , such that

$$\lim_{n \rightarrow \infty} P \left\{ \sup_{n^{-1} \|\mathbf{W}_n\|_2^2 = M_\epsilon^2} \tilde{L}_n(\Psi_n^0 + \gamma_n \mathbf{W}_n; \boldsymbol{\lambda}_n, h_n) < \tilde{L}_n(\Psi_n^0; \boldsymbol{\lambda}_n, h_n) \right\} \geq 1 - \epsilon$$

as needed in (A.1). This completes the proof of Theorem 1.

Proof of Lemma 1

Proof. Recall the partitioning of the true parameter matrix $\Psi_n^0 = (\Psi_{n1}^0, \Psi_{n2}^0)$ such that Ψ_{n2}^0 corresponds to all the zero regression functions, while Ψ_{n1}^0 corresponds to the nonzero regression functions, among other parameters, given in Section 4 of the paper. Also, consider any matrix $\Psi_n = (\Psi_{n1}, \Psi_{n2})$ in a neighbourhood of Ψ_n^0 such that $n^{-1} \|\Psi_n - \Psi_n^0\|_F^2 = O_p((nh_n)^{-1})$, and $\dim(\Psi_{n1}) = \dim(\Psi_{n1}^0)$ and $\dim(\Psi_{n2}) = \dim(\Psi_{n2}^0)$. The existence of such neighbourhood is guaranteed by the result of Theorem 1. We provide the proof in two steps:

Step 1. We first prove that, with probability tending to one, as $n \rightarrow \infty$,

$$\tilde{L}_n((\Psi_{n1}, \Psi_{n2}); \boldsymbol{\lambda}_n, h_n) - \tilde{L}_n((\Psi_{n1}, \mathbf{0}); \boldsymbol{\lambda}_n, h_n) < 0. \tag{A.47}$$

By the definition of the penalized likelihood function in (9) of the paper, we have that

$$\begin{aligned} & \tilde{L}_n((\Psi_{n1}, \Psi_{n2}); \lambda_n, h_n) - \tilde{L}_n((\Psi_{n1}, \mathbf{0}); \lambda_n, h_n) \\ &= [L_n((\Psi_{n1}, \Psi_{n2}); h_n) - L_n((\Psi_{n1}, \mathbf{0}); h_n)] \\ & - [\mathbb{P}_n((\Psi_{n1}, \Psi_{n2}); \lambda_n) - \mathbb{P}_n((\Psi_{n1}, \mathbf{0}); \lambda_n)]. \end{aligned} \quad (\text{A.48})$$

By the definition of the penalty in (10) of the paper, we have

$$\mathbb{P}_n((\Psi_{n1}, \Psi_{n2}); \lambda_n) - \mathbb{P}_n((\Psi_{n1}, \mathbf{0}); \lambda_n) = \sum_{j=1}^C \sum_{l=d_0^j+1}^d p_n(\|\mathbf{b}_{jl}\|_2/\sqrt{n}; \lambda_{nj}). \quad (\text{A.49})$$

On the other hand, using (7) and (8) of the paper, we can write the likelihood difference as

$$\begin{aligned} L_n((\Psi_{n1}, \Psi_{n2}); h_n) - L_n((\Psi_{n1}, \mathbf{0}); h_n) &= \sum_{i=1}^n [\ell_n((\boldsymbol{\psi}_1(u_t), \boldsymbol{\psi}_2(u_t)); h_n) - \ell_n((\boldsymbol{\psi}_1(u_t), \mathbf{0}); h_n)] \\ &= \sum_{t=1}^n \sum_{i=1}^n [l((\boldsymbol{\psi}_1(U_t), \boldsymbol{\psi}_2(U_t)), \mathbf{X}_i, Y_i) - l((\boldsymbol{\psi}_1(U_t), \mathbf{0}), \mathbf{X}_i, Y_i)] K_{h_n}(U_i - U_t), \end{aligned} \quad (\text{A.50})$$

where $l(\cdot)$ is the log of the mixture density as given in (6) of the paper. Note that $\boldsymbol{\psi}(u_t) = (\boldsymbol{\psi}_1(u_t), \boldsymbol{\psi}_2(u_t))$ is the t -th row of the matrix $\Psi_n = (\Psi_{n1}, \Psi_{n2})$. For any fixed $U_t = u \in \mathcal{U}$, we first assess the inner sum in (A.50). We thus consider

$$\Delta_n(u) = \sum_{i=1}^n [l((\boldsymbol{\psi}_1(u), \boldsymbol{\psi}_2(u)), \mathbf{X}_i, Y_i) - l((\boldsymbol{\psi}_1(u), \mathbf{0}), \mathbf{X}_i, Y_i)] K_{h_n}(U_i - u) \quad (\text{A.51})$$

By the mean value theorem, there exists a vector $\boldsymbol{\xi}(u)$ on the segment between $\mathbf{0}$ and $\boldsymbol{\psi}_2(u)$ such that

$$\begin{aligned} & \Delta_n(u) \\ &= \sum_{i=1}^n K_{h_n}(U_i - u) \left[\frac{\partial l((\boldsymbol{\psi}_1(u), \boldsymbol{\xi}(u)), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} \right]^\top \boldsymbol{\psi}_2(u) \\ &= \sum_{i=1}^n K_{h_n}(U_i - u) \left[\frac{\partial l((\boldsymbol{\psi}_1(u), \boldsymbol{\xi}(u)), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} - \frac{\partial l((\boldsymbol{\psi}_1^0(u), \mathbf{0}), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} \right]^\top \boldsymbol{\psi}_2(u) \\ &+ \sum_{i=1}^n K_{h_n}(U_i - u) \left[\frac{\partial l((\boldsymbol{\psi}_1^0(u), \mathbf{0}), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} \right]^\top \boldsymbol{\psi}_2(u) = \Delta_1(u) + \Delta_2(u), \end{aligned} \quad (\text{A.52})$$

where $\Delta_1(u)$ and $\Delta_2(u)$ are respectively the two sums on the right hand side of the above equation. Similar to (A.18), as shown in the proof of Theorem 1 for the order assessment of

(A.5), we have that, for any $u \in \mathcal{U}$,

$$\sum_{i=1}^n \frac{\partial l((\boldsymbol{\psi}_1^0(u), \mathbf{0}), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} K_{h_n}(U_i - u) = O_p\{(nh_n^{-1})^{\frac{1}{2}}\}. \quad (\text{A.53})$$

Thus, since $\dim(\boldsymbol{\psi}_2(u)) < \infty$, using the Cauchy-Schwarz inequality we have

$$\begin{aligned} \Delta_2(u) &= O_p\{(nh_n^{-1})^{\frac{1}{2}}\} \|\boldsymbol{\psi}_2(u)\|_1 \\ &= O_p\{(nh_n^{-1})^{\frac{1}{2}}\} \|\boldsymbol{\psi}_2(u)\|_2 \\ &= O_p\{(nh_n^{-1})^{\frac{1}{2}}\} \left(\sum_{j=1}^C \sum_{l=d_{j+1}^0}^d \beta_{jl}^2(u) \right)^{1/2} \end{aligned} \quad (\text{A.54})$$

which then implies that

$$\sum_{t=1}^n \Delta_2(U_t) = O_p\{(nh_n^{-1})^{\frac{1}{2}}\} \sum_{t=1}^n \|\boldsymbol{\psi}_2(U_t)\|_2$$

or using the Cauchy-Schwarz inequality,

$$\begin{aligned} \sum_{t=1}^n \Delta_2(U_t) &= O_p\left(\frac{n}{\sqrt{h_n}}\right) \left\{ \sum_{t=1}^n \|\boldsymbol{\psi}_2(U_t)\|_2^2 \right\}^{1/2} \\ &= O_p\left(\frac{n^{3/2}}{\sqrt{h_n}}\right) \left\{ \sum_{j=1}^c \sum_{l=d_{j+1}^0}^d \|\mathbf{b}_{jl}\|_2 / \sqrt{n} \right\}. \end{aligned} \quad (\text{A.55})$$

For the first term in (A.52), we have that

$$\begin{aligned} \Delta_1(u) &= \sum_{i=1}^n K_{h_n}(U_i - u) \left[\frac{\partial l((\boldsymbol{\psi}_1(u), \boldsymbol{\xi}(u)), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} - \frac{\partial l((\boldsymbol{\psi}_1(u), \mathbf{0}), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} \right]^\top \boldsymbol{\psi}_2(u) \\ &\quad + \sum_{i=1}^n K_{h_n}(U_i - u) \left[\frac{\partial l((\boldsymbol{\psi}_1(u), \mathbf{0}), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} - \frac{\partial l((\boldsymbol{\psi}_1^0(u), \mathbf{0}), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} \right]^\top \boldsymbol{\psi}_2(u). \end{aligned}$$

By the mean value theorem, there exist vectors $\boldsymbol{\xi}^*(u)$ and $\boldsymbol{\psi}_1^*(u)$ respectively on segments between $\mathbf{0}$ and $\boldsymbol{\xi}(u)$, and $\boldsymbol{\psi}_1(u)$ and $\boldsymbol{\psi}_1^0(u)$ such that, by Condition (RC.3),

$$\begin{aligned} &\sum_{i=1}^n K_{h_n}(U_i - u) \left[\frac{\partial l((\boldsymbol{\psi}_1(u), \boldsymbol{\xi}(u)), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} - \frac{\partial l((\boldsymbol{\psi}_1(u), \mathbf{0}), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} \right] \\ &= \sum_{i=1}^n K_{h_n}(U_i - u) [l''((\boldsymbol{\psi}_1(u), \boldsymbol{\xi}^*(u)), \mathbf{X}_i, Y_i)] \boldsymbol{\xi}(u) = O_p(n) \|\boldsymbol{\xi}(u)\|_2 \end{aligned}$$

and

$$\begin{aligned} & \sum_{i=1}^n K_{h_n}(U_i - u) \left[\frac{\partial l((\boldsymbol{\psi}_1(u), \mathbf{0}), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} - \frac{\partial l((\boldsymbol{\psi}_1^0(u), \mathbf{0}), \mathbf{X}_i, Y_i)}{\partial \boldsymbol{\psi}_2(u)} \right] \\ &= \sum_{i=1}^n K_{h_n}(U_i - u) l''((\boldsymbol{\psi}_1^*(u), \mathbf{0}), \mathbf{X}_i, Y_i) [\boldsymbol{\psi}_1(u) - \boldsymbol{\psi}_1^0(u)] = O_p(n) \|\boldsymbol{\psi}_1(u) - \boldsymbol{\psi}_1^0(u)\|_2. \end{aligned}$$

Thus, together with the above order assessment, and again the Cauchy-Schwarz inequality, we have that

$$\Delta_1(u) = O_p(n) \{ \|\boldsymbol{\xi}(u)\|_2 + \|\boldsymbol{\psi}_1(u) - \boldsymbol{\psi}_1^0(u)\|_2 \} \|\boldsymbol{\psi}_2(u)\|_2. \quad (\text{A.56})$$

This implies that

$$\sum_{t=1}^n \Delta_1(U_t) = O_p(n) \sum_{t=1}^n \{ \|\boldsymbol{\xi}(U_t)\|_2 + \|\boldsymbol{\psi}_1(U_t) - \boldsymbol{\psi}_1^0(U_t)\|_2 \} \|\boldsymbol{\psi}_2(U_t)\|_2. \quad (\text{A.57})$$

By the Cauchy-Schwarz inequality,

$$\begin{aligned} & \sum_{t=1}^n \{ \|\boldsymbol{\xi}(U_t)\|_2 + \|\boldsymbol{\psi}_1(U_t) - \boldsymbol{\psi}_1^0(U_t)\|_2 \} \|\boldsymbol{\psi}_2(U_t)\|_2 \\ & \leq \left\{ \sum_{t=1}^n [\|\boldsymbol{\xi}(U_t)\|_2^2 + \|\boldsymbol{\psi}_1(U_t) - \boldsymbol{\psi}_1^0(U_t)\|_2^2] \right\}^{1/2} \left\{ \sum_{t=1}^n \|\boldsymbol{\psi}_2(U_t)\|_2^2 \right\}^{1/2} \\ & = \{n(nh_n)^{-1}\}^{1/2} \left\{ \sum_{j=1}^C \sum_{l=d_j^0+1}^d \|\mathbf{b}_{jl}\|_2^2 \right\}^{1/2} \leq h_n^{-1/2} \left\{ \sum_{j=1}^C \sum_{l=d_j^0+1}^d \|\mathbf{b}_{jl}\|_2 \right\}. \end{aligned}$$

Combing the above results, for some constant C_1 , we arrive at,

$$\sum_{t=1}^n \Delta_1(U_t) = \frac{C_1 n^{3/2}}{\sqrt{h_n}} \left\{ \sum_{j=1}^C \sum_{l=d_j^0+1}^d \|\mathbf{b}_{jl}\|_2 / \sqrt{n} \right\}. \quad (\text{A.58})$$

By (A.55), (A.58), and (A.50), we have that, for large n ,

$$L_n((\boldsymbol{\Psi}_{n1}, \boldsymbol{\Psi}_{n2}); \lambda_n, h_n) - L_n((\boldsymbol{\Psi}_{n1}, \mathbf{0}); \lambda_n, h_n) = \frac{C_2 n^{3/2}}{\sqrt{h_n}} \left\{ \sum_{j=1}^C \sum_{l=d_j^0+1}^d \|\mathbf{b}_{jl}\|_2 / \sqrt{n} \right\} \quad (\text{A.59})$$

for some constant C_2 . Finally, by using (A.49) and (A.59) in (A.48), for large n ,

$$\begin{aligned} & \tilde{L}_n((\Psi_{n1}, \Psi_{n2}); \lambda_n, h_n) - \tilde{L}_n((\Psi_{n1}, \mathbf{0}); \lambda_n, h_n) \\ &= \frac{C_2 n^{3/2}}{\sqrt{h_n}} \left\{ \sum_{j=1}^C \sum_{l=d_j^0+1}^d \|\mathbf{b}_{jl}\|_2 / \sqrt{n} \right\} - \sum_{j=1}^C \sum_{l=d_j^0+1}^d p_n(\|\mathbf{b}_{jl}\|_2 / \sqrt{n}; \lambda_n) \\ &= \sum_{j=1}^C \sum_{l=d_j^0+1}^d \left\{ \frac{C_2 n^{3/2}}{\sqrt{h_n}} \|\mathbf{b}_{jl}\|_2 / \sqrt{n} - p_n(\|\mathbf{b}_{jl}\|_2 / \sqrt{n}; \lambda_n) \right\} < 0. \end{aligned}$$

The last inequality is due to Condition (PC.3) on the penalty function p_n . This completes the proof of (A.47).

Step 2. Consider the penalized local log-likelihood function $\tilde{L}_n((\Psi_{n1}, \mathbf{0}); \lambda_n, h_n)$, which is only a function of Ψ_{n1} . Let $(\hat{\Psi}_{n1}, \mathbf{0})$ be its maximizer. Then for any $\Psi_n = (\Psi_{n1}, \Psi_{n2})$ with the property (A.47) shown in **Step 1**, we have that,

$$\begin{aligned} & L_n((\Psi_{n1}, \Psi_{n2}); \lambda_n, h_n) - L_n((\hat{\Psi}_{n1}, \mathbf{0}); \lambda_n, h_n) \\ &= [L_n((\Psi_{n1}, \Psi_{n2}); \lambda_n, h_n) - L_n((\Psi_{n1}, \mathbf{0}); \lambda_n, h_n)] \\ &\quad - [L_n((\hat{\Psi}_{n1}, \mathbf{0}); \lambda_n, h_n) - L_n((\Psi_{n1}, \mathbf{0}); \lambda_n, h_n)] \\ &\leq L_n((\Psi_{n1}, \Psi_{n2}); \lambda_n, h_n) - L_n((\hat{\Psi}_{n1}, \mathbf{0}); \lambda_n, h_n) < 0, \end{aligned}$$

where the last two inequalities are due to the definition of $(\hat{\Psi}_{n1}, \mathbf{0})$ and the result of **Step 1**. Hence, in the $\sqrt{nh_n}$ -neighbourhood of Ψ_n^0 guaranteed by Theorem 1, the maximizer of the penalized local log-likelihood $L_n((\Psi_{n1}, \Psi_{n2}); \lambda_n, h_n)$ satisfies $\hat{\Psi}_{n2} = \mathbf{0}$, with probability tending to one, as $n \rightarrow \infty$. This completes the proof of Lemma 1.

Proof of Theorem 2

Proof. (i) Recall the partitioning of the true parameter matrix $\Psi_n^0 = (\Psi_{n1}^0, \Psi_{n2}^0)$ such that Ψ_{n2}^0 corresponds to all the zero regression functions, while Ψ_{n1}^0 corresponds to the nonzero regression functions, among other parameters, given in Section 4 of the paper. Consider any matrix $\Psi_n = (\Psi_{n1}, \Psi_{n2})$ in a neighbourhood of Ψ_n^0 such that $n^{-1} \|\Psi_n - \Psi_n^0\|_F^2 = O_p((nh_n)^{-1})$, and $\dim(\Psi_{n1}) = \dim(\Psi_{n1}^0)$ and $\dim(\Psi_{n2}) = \dim(\Psi_{n2}^0)$. Note that such choice is guaranteed

by Theorem 1. Also, for the t -th row of this matrix we have that $\|\boldsymbol{\psi}(U_t) - \boldsymbol{\psi}^0(U_t)\|_2 = O_p((nh_n)^{-1/2})$, $t = 1, \dots, n$. On the other hand, for any $u \in \mathcal{U}$, let u^* be its nearest neighbourhood among the observed index values U_1, \dots, U_n , that is $u^* = \operatorname{argmin}_{\tilde{u} \in \{U_1, U_2, \dots, U_n\}} |u - \tilde{u}|$. Under Condition (RC.2), we have that $\|\boldsymbol{\psi}^0(u) - \boldsymbol{\psi}^0(u^*)\|_2 = O_p(\log n/n)$, (Janson, 1987). Using the triangle inequality, we then have that $\|\boldsymbol{\psi}(u) - \boldsymbol{\psi}^0(u)\|_2 \leq \|\boldsymbol{\psi}(u) - \boldsymbol{\psi}(u^*)\|_2 + \|\boldsymbol{\psi}(u^*) - \boldsymbol{\psi}^0(u^*)\|_2 + \|\boldsymbol{\psi}^0(u^*) - \boldsymbol{\psi}^0(u)\|_2$. Thus, for any $u \in \mathcal{U}$, we have that $\|\boldsymbol{\psi}(u) - \boldsymbol{\psi}^0(u)\|_2 = O_p((nh_n)^{-1/2})$. We provide the proof in two steps as follows.

Step 1: For any $u \in \mathcal{U}$, using the partitioning $\boldsymbol{\psi}(u) = (\boldsymbol{\psi}_1(u), \boldsymbol{\psi}_2(u))$ for any choice $\boldsymbol{\psi}(u)$ in the above neighbourhood, consider the local log-likelihood given in (7) of the paper,

$$\ell_n((\boldsymbol{\psi}_1(u), \boldsymbol{\psi}_2(u)); h_n) = \sum_{i=1}^n l((\boldsymbol{\psi}_1(u), \boldsymbol{\psi}_2(u)), \mathbf{X}_i, Y_i) K_{h_n}(U_i - u).$$

We now denote the following two penalized local log-likelihoods,

$$\begin{aligned} \tilde{\ell}_n(\boldsymbol{\psi}(u); h_n, \boldsymbol{\lambda}_n) &= \ell_n((\boldsymbol{\psi}_1(u), \boldsymbol{\psi}_2(u)); h_n) - \frac{1}{n} \sum_{j=1}^C \sum_{l=1}^d p_n(\|\mathbf{b}_{jl}\|_2/\sqrt{n}; \lambda_{nj}), \\ \tilde{\ell}_n(\boldsymbol{\psi}_1(u); h_n, \boldsymbol{\lambda}_n) &= \ell_n((\boldsymbol{\psi}_1(u), \mathbf{0}); h_n) - \frac{1}{n} \sum_{j=1}^C \sum_{l=1}^{d_j^0} p_n(\|\mathbf{b}_{jl}\|_2/\sqrt{n}; \lambda_{nj}). \end{aligned}$$

In what follows, for any $u \in \mathcal{U}$, we assess the order of the following difference for large n ,

$$\begin{aligned} D_n(u) &= \tilde{\ell}_n(\boldsymbol{\psi}(u); h_n, \boldsymbol{\lambda}_n) - \tilde{\ell}_n(\boldsymbol{\psi}_1(u); h_n, \boldsymbol{\lambda}_n) \\ &= [\ell_n((\boldsymbol{\psi}_1(u), \boldsymbol{\psi}_2(u)); h_n) - \ell_n((\boldsymbol{\psi}_1(u), \mathbf{0}); h_n)] \\ &\quad - \frac{1}{n} \sum_{j=1}^C \sum_{l=d_j^0+1}^d p_n(\|\mathbf{b}_{jl}\|_2/\sqrt{n}; \lambda_{nj}) \\ &= \Delta_n(u) - \frac{1}{n} \sum_{j=1}^C \sum_{l=d_j^0+1}^d p_n(\|\mathbf{b}_{jl}\|_2/\sqrt{n}; \lambda_{nj}), \end{aligned}$$

where $\Delta_n(u)$ is the difference in the local log-likelihood which is also used in the proof of Lemma 1, given in (A.51). Recall the following representation of $\Delta_n(u)$ given in (A.52),

$$\Delta_n(u) = \Delta_1(u) + \Delta_2(u).$$

From (A.54), we have that (up to a constant), for large enough n

$$|\Delta_2(u)| \leq \sqrt{\frac{n}{h_n}} \sum_{j=1}^C \sum_{l=d_{j+1}^0}^d |\beta_{jl}(u)|.$$

Under Condition (RC.2), we have that $\|\boldsymbol{\psi}^0(u) - \boldsymbol{\psi}^0(u^*)\|_2 = O_p(\frac{\log n}{n})$. Using (A.56), for large n , we also have that (up to a constant), for any $u \in \mathcal{U}$,

$$|\Delta_1(u)| \leq \sqrt{\frac{n}{h_n}} \sum_{j=1}^C \sum_{l=d_{j+1}^0}^d |\beta_{jl}(u)|.$$

Therefore, for any $u \in \mathcal{U}$, for large n we have

$$D_n(u) \leq \sqrt{\frac{n}{h_n}} \sum_{j=1}^C \sum_{l=d_{j+1}^0}^d \left\{ |\beta_{jl}(u)| - \frac{\sqrt{h_n} p_n(\|\mathbf{b}_{jl}\|_2/\sqrt{n}; \lambda_{nj})}{n^{3/2}} \right\}.$$

Also, for any $u \in \mathcal{U}$ and its corresponding u^* as defined above, for large n we have that

$$\begin{aligned} |\beta_{jl}(u)| &\leq |\beta_{jl}(u) - \beta_{jl}(u^*)| + |\beta_{jl}(u^*)| \\ &= O_p(\log n/n) + \|\mathbf{b}_{jl}\|_2/\sqrt{n} + O_p(\{nh_n\}^{-1/2}) \\ &= \|\mathbf{b}_{jl}\|_2/\sqrt{n} + o_p(1). \end{aligned}$$

Therefore, for large n we have that by the Condition (PC.3) on the penalty, $D_n(u) < 0$, uniformly in $u \in \mathcal{U}$.

Step 2: For any $u \in \mathcal{U}$, let $\widehat{\boldsymbol{\psi}}_{n1}(u)$ be the maximizer of $\tilde{\ell}_n(\boldsymbol{\psi}_1(u); h_n, \boldsymbol{\lambda}_n)$ which is considered as only a function of $\boldsymbol{\psi}_1(u)$. Then, for any $u \in \mathcal{U}$,

$$\begin{aligned} &\tilde{\ell}_n(\boldsymbol{\psi}(u); h_n, \boldsymbol{\lambda}_n) - \tilde{\ell}_n(\widehat{\boldsymbol{\psi}}_{n1}(u); h_n, \boldsymbol{\lambda}_n) \\ &= \left\{ \tilde{\ell}_n(\boldsymbol{\psi}(u); h_n, \boldsymbol{\lambda}_n) - \tilde{\ell}_n(\boldsymbol{\psi}_1(u); h_n, \boldsymbol{\lambda}_n) \right\} \\ &\quad - \left\{ \tilde{\ell}_n(\widehat{\boldsymbol{\psi}}_{n1}(u); h_n, \boldsymbol{\lambda}_n) - \tilde{\ell}_n(\boldsymbol{\psi}_1(u); h_n, \boldsymbol{\lambda}_n) \right\} < 0. \end{aligned}$$

This difference is indeed negative uniformly over $u \in \mathcal{U}$. This implies that with probability tending to one, as $n \rightarrow \infty$, $\widehat{\boldsymbol{\psi}}_{n2}(u) = \mathbf{0}$, uniformly in u . Similar to the result of Theorem 1, $\widehat{\boldsymbol{\psi}}_{n1}(u)$ is a consistent estimator of $\boldsymbol{\psi}_1^0(u)$, for any $u \in \mathcal{U}$.

(ii) For any $u \in \mathcal{U}$, recall the oracle estimator $\widehat{\boldsymbol{\psi}}_{n1,orc}(u)$ which is the MLLE of $\boldsymbol{\psi}_1^0(u)$, having known $\boldsymbol{\psi}_2^0(u) = \mathbf{0}$ a priori. This estimator is the maximizer of the local log-likelihood

$\ell(\boldsymbol{\psi}(u); h_n)$ in (7) of the paper, and it satisfies the estimating equation

$$\ell'_n(\widehat{\boldsymbol{\psi}}_{n1,orc}(u); h_n) = \mathbf{0}.$$

Using a Taylor’s expansion around $\boldsymbol{\psi}_1^0(u)$, we arrive at

$$\ell'_n(\boldsymbol{\psi}_1^0(u); h_n) + \ell''_n(\tilde{\boldsymbol{\psi}}_1^0(u); h_n)(\widehat{\boldsymbol{\psi}}_{n1,orc}(u) - \boldsymbol{\psi}_1^0(u)) = \mathbf{0},$$

where $\tilde{\boldsymbol{\psi}}_1^0(u)$ lies on a segment between $\boldsymbol{\psi}_1^0(u)$ and $\widehat{\boldsymbol{\psi}}_{n1,orc}(u)$, for any $u \in \mathcal{U}$. By Proposition 1 of the paper, since $\widehat{\boldsymbol{\psi}}_{n1,orc}(u)$ is a consistent estimator of $\boldsymbol{\psi}_1^0(u)$, so is $\tilde{\boldsymbol{\psi}}_1^0(u)$. Thus, for large n ,

$$\ell'_n(\boldsymbol{\psi}_1^0(u); h_n) + \{\ell''_n(\boldsymbol{\psi}_1^0(u); h_n) + o_p(n)\}(\widehat{\boldsymbol{\psi}}_{n1,orc}(u) - \boldsymbol{\psi}_1^0(u)) = \mathbf{0}. \quad (\text{A.60})$$

Similarly, the MPLLE $\widehat{\boldsymbol{\psi}}_{n1}(u)$ also satisfies the estimation equation

$$\tilde{\ell}'_n(\widehat{\boldsymbol{\psi}}_{n1}(u); h_n, \boldsymbol{\lambda}_n) = \mathbf{0}$$

or equivalently,

$$\ell'_n(\widehat{\boldsymbol{\psi}}_{n1}(u); h_n) - \mathbf{P}'_n(\widehat{\boldsymbol{\Psi}}_{n1}; \boldsymbol{\lambda}_n) = \mathbf{0},$$

where

$$\mathbf{P}_n(\boldsymbol{\Psi}_{n1}; \boldsymbol{\lambda}_n) = \frac{1}{n} \sum_{j=1}^C \sum_{l=1}^{d_j^0} p_n(\|\mathbf{b}_{jl}\|_2 / \sqrt{n}; \lambda_{nj}) = \mathbb{P}_n(\boldsymbol{\Psi}_{n1}; \boldsymbol{\lambda}_n) / n. \quad (\text{A.61})$$

Using a Taylor’s expansion around $\boldsymbol{\psi}_1^0(u)$, and consistency of $\widehat{\boldsymbol{\psi}}_{n1}(u)$, for large n ,

$$\begin{aligned} & \{\ell'_n(\boldsymbol{\psi}_1^0(u); h_n) - [\mathbf{P}'_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n) + o_p(n)]\} \\ & + \{\ell''_n(\boldsymbol{\psi}_1^0(u); h_n) - \mathbf{P}''_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n) + o_p(n)\}(\widehat{\boldsymbol{\psi}}_{n1}(u) - \boldsymbol{\psi}_1^0(u)) \\ & = \mathbf{0}, \end{aligned} \quad (\text{A.62})$$

where \mathbf{P}'_n and \mathbf{P}''_n are the gradient and Hessian of the penalty \mathbf{P}_n in (A.61) with respect to the parameter vector $\boldsymbol{\psi}_1(u)$, and $\boldsymbol{\theta}_1^0 = \{\theta_{jl}^0 = \sqrt{E_0\{\beta_{jl}^0(U)\}^2} : 1 \leq j \leq C, 1 \leq l \leq d_j^0\}$.

By (A.60) and (A.62), we have that

$$\begin{aligned}
& \widehat{\boldsymbol{\psi}}_{n1}(u) - \widehat{\boldsymbol{\psi}}_{n1,orc}(u) \\
&= \left\{ \ell_n''(\boldsymbol{\psi}_1^0(u); h_n) + o_p(n) \right\}^{-1} \ell_n'(\boldsymbol{\psi}_1^0(u); h_n) \\
&\quad - \left\{ \ell_n''(\boldsymbol{\psi}_1^0(u); h_n) - \mathbf{P}_n''(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n) + o_p(n) \right\}^{-1} \left\{ \ell_n'(\boldsymbol{\psi}_1^0(u); h_n) - \mathbf{P}_n'(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n) \right\} \\
&= \left\{ \left\{ \ell_n''(\boldsymbol{\psi}_1^0(u); h_n)/n + o_p(1) \right\}^{-1} - \left\{ \ell_n''(\boldsymbol{\psi}_1^0(u); h_n)/n - \mathbf{P}_n''(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)/n + o_p(1) \right\}^{-1} \right\} \\
&\quad \times \left\{ \frac{\ell_n'(\boldsymbol{\psi}_1^0(u); h_n)}{n} \right\} \\
&\quad + \left\{ \ell_n''(\boldsymbol{\psi}_1^0(u); h_n)/n - \mathbf{P}_n''(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)/n + o_p(1) \right\}^{-1} \left\{ \frac{\mathbf{P}_n'(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)}{n} + o_p(1) \right\} \\
&= \mathbf{V}_{1n}(u) + \mathbf{V}_{2n}(u),
\end{aligned}$$

where $\mathbf{V}_{1n}(u)$ and $\mathbf{V}_{2n}(u)$ are respectively the two vectors on the right hand side of the above equation. By (A.25), for any $u \in \mathcal{U}$ and for large n ,

$$\ell_n''(\boldsymbol{\psi}_1^0(u); h_n)/n = -m(u)\mathbf{I}(u)(1 + o_p(1)). \quad (\text{A.63})$$

Thus, due to the finiteness condition of $m(u)\mathbf{I}(u)$ and since $r_{2n} = o(1)$ by Condition (PC.2) on the penalty, we have (in a matrix form)

$$\left\{ \ell_n''(\boldsymbol{\psi}_1^0(u); h_n)/n + o_p(1) \right\}^{-1} - \left\{ \ell_n''(\boldsymbol{\psi}_1^0(u); h_n)/n - \mathbf{P}_n''(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)/n + o_p(1) \right\}^{-1} = o_p(1). \quad (\text{A.64})$$

Also, by (A.14) we have

$$E \left\{ \frac{\ell_n'(\boldsymbol{\psi}_1^0(u); h_n)}{n} \right\} = \frac{1}{2} \mathcal{K}_2 m(u) \left\{ \Delta''(u; u) + 2\Delta'(u; u)m'(u)/m(u) \right\} h_n^2 \{1 + o(1)\} < \infty \quad (\text{A.65})$$

for any $u \in \mathcal{U}$, and $\mathcal{K}_2 = \int_{\mathcal{U}} t^2 K(t) dt < \infty$. Using a similar approach to show (A.26) and the centralizing technique as in (A.36), and since $nh_n^5 = O(1)$, we have that

$$\sup_{u \in \mathcal{U}} \left\| \ell_n'(\boldsymbol{\psi}_1^0(u); h_n)/n \right\|_2 = O_p \left\{ \left(\frac{\log(\frac{1}{h_n})}{nh_n} \right)^{\frac{1}{2}} \right\} + O_p(h_n^2) = O_p \left\{ \left(\frac{\log(\frac{1}{h_n})}{nh_n} \right)^{\frac{1}{2}} \right\}. \quad (\text{A.66})$$

Thus, by (A.64) and (A.66),

$$\sup_{u \in \mathcal{U}} \left\| \mathbf{V}_{1n}(u) \right\|_2 = o_p \left\{ \left(\frac{\log(\frac{1}{h_n})}{nh_n} \right)^{\frac{1}{2}} \right\} = o_p(\xi_n). \quad (\text{A.67})$$

Also, by the positive definiteness of $m(u)\mathbf{I}(u)$ and condition $r_{2n} = o(1)$, we have that

$$\sup_{u \in \mathcal{U}} \|\mathbf{V}_{2n}(u)\|_2 = O_p \left\{ r_{1n} / \sqrt{nh_n} \right\}. \quad (\text{A.68})$$

Thus, using (A.67) and (A.68) we have

$$\begin{aligned} & \{(1 + r_{1n})\xi_n\}^{-1} \sup_{u \in \mathcal{U}} \|\widehat{\boldsymbol{\psi}}_{n1}(u) - \widehat{\boldsymbol{\psi}}_{n1,orc}(u)\|_2 \\ & \leq \{(1 + r_{1n})\xi_n\}^{-1} \left\{ \sup_{u \in \mathcal{U}} \|\mathbf{V}_{1n}(u)\|_2 + \sup_{u \in \mathcal{U}} \|\mathbf{V}_{2n}(u)\|_2 \right\} \\ & = \{(1 + r_{1n})\}^{-1} o_p(1) + \frac{r_{1n}}{1 + r_{1n}} \sqrt{\frac{1}{\log h_n^{-1}}} = o_p(1). \end{aligned}$$

This completes the proof.

(iii) From (A.62) and using (A.63), and the relation (A.61) we have

$$\left\{ m(u)\mathbf{I}(u) + \frac{\mathbb{P}'_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)}{n^2} + o_p(1) \right\} \left(\widehat{\boldsymbol{\psi}}_{n1}(u) - \boldsymbol{\psi}_1^0(u) \right) + \frac{\mathbb{P}_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)}{n^2} = \frac{\ell'_n(\boldsymbol{\psi}_1^0(u); h_n)}{n}. \quad (\text{A.69})$$

Using the second-order moment calculations in the proof of Theorem 2, we have

$$\text{Var} \left\{ \frac{\ell'_n(\boldsymbol{\psi}_1^0(u); h_n)}{n} \right\} = \frac{\mathcal{V}_0}{nh_n} \mathbf{I}(u) m(u) (1 + o_p(1)) < \infty, \quad (\text{A.70})$$

where $\mathcal{V}_0 = \int_{\mathcal{U}} K^2(t) dt < \infty$.

Using (A.65) and (A.70), we standardize both sides of the equation (A.69). Hence, by the multivariate central limit theorem we have that

$$\sqrt{nh_n} \left\{ \left[\mathbf{I}(u) + m^{-1}(u) \frac{\mathbb{P}''_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)}{n^2} \right] \left(\widehat{\boldsymbol{\psi}}_{n1}(u) - \boldsymbol{\psi}_1^0(u) \right) - \mathbf{B}_n(u) \right\} \sim N \left(\mathbf{0}, \mathcal{V}_0 m^{-1}(u) \mathbf{I}(u) \right),$$

where

$$\mathbf{B}_n(u) = \frac{1}{2} \mathcal{K}_2 \left\{ \Delta''(u; u) + 2\Delta'(u; u) m'(u) / m(u) \right\} h_n^2 - m^{-1}(u) \frac{\mathbb{P}'_n(\boldsymbol{\theta}_1^0; \boldsymbol{\lambda}_n)}{n^2} + o_p(h_n^2).$$

Web Appendix E. Remarks

REMARK 1: The sparsity assumption in Equation (5) of the paper is indeed the most common assumption in the variable selection literature. The implication is that the underlying data-generating process is a simple sparse model and is beneficial for interpretation

purposes, specifically in mixture regressions when the number of covariates is relatively large. Together with (5) and under Conditions (PC.2) and (PC.3) on the penalty function p_n in Web Appendix B, we achieve the selection consistency property as stated in Theorem 2(i). Essentially, Conditions (PC.2) and (PC.3) imply that only those coefficients that satisfy $E\{\beta_{jl}(U)\} > \lambda_n$, where $\lambda_n \rightarrow 0$ as $n \rightarrow \infty$, are detectable by the regularization methods and will be estimated as non-zero, while the ones below the threshold λ_n , so-called weak signals, will most likely be estimated as zero. In a sense, the aforementioned conditions are similar to the ones considered by Wei et al. (2011) with η_1 replaced by λ_n . Wei et al. (2011) refer to (5) as the narrow sparsity condition ($\eta_1 = 0$) under which they also show that, for example, AdpLASSO has the selection consistency property. By changing (5) to the one used in Wei et al. (2011) with $\eta_1 > 0$, it turns out that regularization methods only achieve certain estimation error bounds but not really selection consistency as those weak-signal regression parameters most likely will be estimated as zero. Hence, since our main result focuses on the variable selection consistency property, we have decided to keep the current definition of sparsity (5) in the FM-VCR models. Nevertheless, it is worth noting that Fang et al. (2021) proposed a two-step procedure based on both variable selection and ridge regression estimators that were shown to be capable of detecting weak signals and providing an estimation of both strong and weak signals.

REMARK 2: The asymptotic normality result in Theorem 2(iii) is obtained as if the true sparse structure of the model is known in advance, i.e., an oracle's perspective. In general, its use in practice to perform likelihood-ratio type inference is reserved (referred to as naive inference) as the true sparse structure of the model is not known in advance and it is estimated by the penalization method. As a result, in the partitioning $\hat{\boldsymbol{\psi}}_n(u) = (\hat{\boldsymbol{\psi}}_{n1}(u), \hat{\boldsymbol{\psi}}_{n2}(u))$ introduced before Theorem 2 of the paper (oracle perspective), in practice due to the variable selection stage the dimension of the sub-vector $\hat{\boldsymbol{\psi}}_{n1}(u)$ is indeed random and may not be equal

to the dimension of true non-zero vector $\boldsymbol{\psi}_1^0(u)$, and hence asymptotically normal distribution may be distorted. The extra variability due to the variable selection needs to be taken into account for further inference and is a part of the general topic of post-selection inference (PoSI, Berk et al. (2013)), which sees a surge of research in recent years for (generalized) linear regression models (see Zhang et al. (2022) for a recent survey), although little for the FMR and FM-VCR models.

REMARK 3: Recall the penalty p_n in (10) of the paper. If the AdpLASSO is used, we have

$$p_n(\|\mathbf{b}_{jl}\|_2/\sqrt{n}; \lambda_n)/n^2 = \lambda_n w_{jl} \|\mathbf{b}_{jl}\|_2/\sqrt{n}, \quad j = 1, \dots, C, \quad l = 1, \dots, d.$$

In practice, one needs to specify the weights w_{jl} . As in Zou (2006), we use $w_{jl} = \{\|\check{\mathbf{b}}_{jl}\|_2/\sqrt{n}\}^{-\gamma}$, for some $\gamma > 0$, where $\|\check{\mathbf{b}}_{jl}\|_2^2 = \sum_{t=1}^n \check{\beta}_{jl}^2(u_t)$ and $\check{\beta}_{jl}(u_t)$'s are the consistent MLE discussed in Section 3 of the paper. Thus, as $n \rightarrow \infty$, the condition $\lambda_n \max_{1 \leq j \leq C, 1 \leq l \leq d_j^0} w_{jl} = o_p(\{nh_n\}^{-1/2})$ in Theorem 1 becomes $\sqrt{nh_n} \lambda_n \rightarrow 0$, and the sparsity condition $\sqrt{nh_n} \lambda_n \min_{1 \leq j \leq C, d_j^0+1 \leq l \leq d} w_{jl} \rightarrow \infty$ in Theorem 2 becomes $(nh_n)^{(\gamma+1)/2} \lambda_n \rightarrow \infty$.

REMARK 4: For the AdpLASSO, the difference $D_{n,II}(\mathbf{W}_n)$ in (A.41) has a simpler upper bound as follows. By the definition of the penalty, we have that,

$$\begin{aligned} D_{n,II}(\mathbf{W}_n) &\leq h_n n^{-1} \sum_{j=1}^C \sum_{l=1}^{d_j^0} \{p_n(\|\mathbf{b}_{jl}^0 + \gamma_n \mathbf{v}_{jl}\|/\sqrt{n}; \lambda_{nj}) - p_n(\|\mathbf{b}_{jl}^0\|/\sqrt{n}; \lambda_{nj})\} \\ &\leq h_n n^{-1} \sum_{j=1}^C \sum_{j=1}^{d_j^0} n^2 \gamma_n \lambda_{nj} w_{jl} \|\mathbf{v}_{jl}\|/\sqrt{n} \\ &\leq M_\epsilon \{nh_n\}^{1/2} \max_{1 \leq j \leq C, 1 \leq l \leq d_j^0} \{\lambda_{nj} w_{jl}\} \end{aligned}$$

which is dominated by the likelihood difference $D_{n,I}(\mathbf{W}_n)$ in (A.2) under the condition

$$\{nh_n\}^{1/2} \max_{1 \leq j \leq C, 1 \leq l \leq d_j^0} \{\lambda_{nj} w_{jl}\} = o_p(1).$$

REMARK 5: For the AdpLASSO, the penalty difference in (A.49) becomes,

$$\mathbb{P}_n((\boldsymbol{\Psi}_{n1}, \boldsymbol{\Psi}_{n2}); \boldsymbol{\lambda}_n) - \mathbb{P}_n((\boldsymbol{\Psi}_{n1}, \mathbf{0}); \boldsymbol{\lambda}_n) = \sum_{j=1}^C \sum_{j=d_j^0+1}^d n^2 \lambda_{nj} w_{jl} (\|\mathbf{b}_{jl}\|/\sqrt{n}).$$

Thus, the difference in (A.48) is bounded by

$$\begin{aligned} & \tilde{L}_n((\Psi_{n1}, \Psi_{n2}); \lambda_n, h_n) - \tilde{L}_n((\Psi_{n1}, \mathbf{0}); \lambda_n, h_n) \\ & \leq \sum_{j=1}^C \sum_{l=d_j^0+1}^d \left\{ \frac{C_2 n^{3/2}}{\sqrt{h_n}} \|\mathbf{b}_{jl}\|_2 / \sqrt{n} - n^2 \min_{1 \leq j \leq C, d_j^0+1 \leq l \leq d} \{\lambda_{nj} w_{jl}\} (\|\mathbf{b}_{jl}\| / \sqrt{n}) \right\} < 0 \end{aligned}$$

under the condition $\sqrt{nh_n} \min_{1 \leq j \leq C, d_j^0+1 \leq l \leq d} \{\lambda_{nj} w_{jl}\} \rightarrow \infty$, as $n \rightarrow \infty$.

Web Appendix F. Details of the numerical algorithm and implementation

Web Appendix F.1 Local quadratic approximation of the penalty functions

The Local Quadratic Approximation (LQA) of a penalty function is given as follows.

$$\begin{aligned} & \bar{\mathbb{P}}_n(\Psi; \lambda) \\ & = \sum_{j=1}^C \sum_{l=1}^d \left\{ p_n \left(\|\mathbf{b}_{jl}^{(m)}\|_2 / \sqrt{n}; \lambda_j \right) \right. \\ & \quad \left. + \frac{p'_n \left(\|\mathbf{b}_{jl}^{(m)}\|_2 / \sqrt{n}; \lambda_j \right)}{2 \|\mathbf{b}_{jl}^{(m)}\|_2 / \sqrt{n}} \left(\mathbf{b}_{jl}^\top \mathbf{b}_{jl} / n - \mathbf{b}_{jl}^{(m)\top} \mathbf{b}_{jl}^{(m)} / n \right) \right\}. \end{aligned} \quad (\text{A.71})$$

For each (j, l) the derivative is given as

$$\frac{\partial \bar{\mathbb{P}}_n(\Psi; \lambda)}{\partial \mathbf{b}_{jl}} = \frac{p'_n \left(\|\mathbf{b}_{jl}^{(m)}\|_2 / \sqrt{n}; \lambda_j \right)}{\|\mathbf{b}_{jl}^{(m)}\|_2 / \sqrt{n}} \times \frac{\mathbf{b}_{jl}}{n}.$$

Specifically, for the well-known penalty functions that are used in the main paper, we have the following expressions:

- LASSO:

$$\frac{\partial \bar{\mathbb{P}}_n(\Psi; \lambda)}{\partial \mathbf{b}_{jl}} = \frac{n^{3/2} \lambda_j}{\|\mathbf{b}_{jl}^{(m)}\|_2} \mathbf{b}_{jl}.$$

- AdpLASSO: Given the weights w_{jl} discussed in Remark 2 of the paper, which are based in the MLLE,

$$\frac{\partial \bar{\mathbb{P}}_n(\Psi; \lambda)}{\partial \mathbf{b}_{jl}} = \frac{n^{3/2} \lambda_j w_{jl}}{\|\mathbf{b}_{jl}^{(m)}\|_2} \mathbf{b}_{jl}.$$

- MCP:

$$\frac{\partial \bar{\mathbb{P}}_n(\Psi; \lambda)}{\partial \mathbf{b}_{jl}} = \begin{cases} \left[\frac{n^{3/2}\lambda_j}{\|\mathbf{b}_{jl}^{(m)}\|_2} - \frac{n}{\gamma} \right] \mathbf{b}_{jl}, & \|\mathbf{b}_{jl}^{(m)}\|_2 \leq \sqrt{n\gamma}\lambda_j \\ \mathbf{0}, & \|\mathbf{b}_{jl}^{(m)}\|_2 > \sqrt{n\gamma}\lambda_j \end{cases}$$

- SCAD:

$$= \begin{cases} \frac{n^{3/2}\lambda_j}{\|\mathbf{b}_{jl}^{(m)}\|_2} \mathbf{b}_{jl} & \|\mathbf{b}_{jl}^{(m)}\|_2 \leq \sqrt{n}\lambda_j \\ \frac{n}{\gamma-1} \left(\frac{\sqrt{n\gamma}\lambda_j}{\|\mathbf{b}_{jl}^{(m)}\|_2} - 1 \right) \mathbf{b}_{jl}, & \sqrt{n}\lambda_j < \|\mathbf{b}_{jl}^{(m)}\|_2 \leq \sqrt{n\gamma}\lambda_j \\ \mathbf{0}, & \|\mathbf{b}_{jl}^{(m)}\|_2 > \sqrt{n\gamma}\lambda_j \end{cases} \quad (\text{A.72})$$

Web Appendix F.2 The EM algorithm for the MLE and MPLLE in Gaussian FM-VCR

For a Gaussian FM-VCR, at any $u \in \mathcal{U}$, the local-kernel log-likelihood in (7) of the paper is given by

$$\ell_n(\boldsymbol{\psi}(u); h) = \sum_{i=1}^n \log \left\{ \sum_{j=1}^C \pi_j N(y_i; \mathbf{x}_i^\top \boldsymbol{\beta}_j, \sigma_j^2) \right\} K_h(u_i - u), \quad (\text{A.73})$$

where $N(y; \mu, \sigma^2)$ is the pdf of a Gaussian distribution with mean μ and σ^2 .

The locally constant vector of parameters is $\boldsymbol{\psi}(u) = (\boldsymbol{\pi}^\top, \boldsymbol{\phi}^\top, \boldsymbol{\beta}_1^\top, \dots, \boldsymbol{\beta}_C^\top)$, where $\boldsymbol{\pi}^\top = (\pi_1, \dots, \pi_{C-1})$, $\boldsymbol{\phi}^\top = (\sigma_1^2, \dots, \sigma_C^2)$ and $\boldsymbol{\beta}_j^\top = (\beta_{j1}, \dots, \beta_{jd})$ for $j = 1, \dots, C$. Note that the entries of the vector $\boldsymbol{\psi}(u)$ are local constant approximations of the functions $\beta_{jl}(u)$, $l = 1, \dots, d$, $\pi_j(u)$, and $\sigma_j^2(u)$, for all $j = 1, \dots, C$. These entries clearly depend on u , and for simplicity, we suppress u in the notation but keep it for $\boldsymbol{\psi}(u)$.

Using (A.73), the corresponding (total) local-kernel log-likelihood in (8) of the paper is

given by

$$\begin{aligned} L_n(\Psi; h) &= \sum_{t=1}^n \ell_n(\boldsymbol{\psi}(u_t); h) \\ &= \sum_{t,i=1}^n \log \left\{ \sum_{j=1}^C \pi_{j,t} N(y_i; \mathbf{x}_i^\top \boldsymbol{\beta}_{j,t}, \sigma_{j,t}^2) \right\} K_h(u_i - u_t), \end{aligned} \quad (\text{A.74})$$

where Ψ is the $(n \times \{C(d+2) - 1\})$ -dimensional matrix of the locally constant vectors $\boldsymbol{\psi}(u_i), i = 1, \dots, n$.

Web Appendix F.2.1 *The EM algorithm for the MLLE.* We view the observed data $\{(u_i, \mathbf{x}_i, y_i) : i = 1, \dots, n\}$ as incomplete, and introduce the unobserved Bernoulli random variables Z_{ij} to represent the membership of the i -th observation to the j -th component of the mixture model, $\forall j = 1, \dots, C$, such that $\Pr(Z_{ij} = 1 | u_i, \mathbf{x}_i, y_i) = \pi_j(u_i)$. The complete data consist of $\{(u_i, \mathbf{x}_i, y_i, \mathbf{z}_i) : i = 1, \dots, n\}$, where $\mathbf{z}_i = (z_{i1}, \dots, z_{iC})^\top$.

At any arbitrary point $u \in \mathcal{U}$, the complete local log-likelihood function is given by

$$\begin{aligned} \ell_n^c(\boldsymbol{\psi}(u); h) &= \sum_{i=1}^n \sum_{j=1}^C Z_{ij} \{ \log \pi_j + \log N(y_i; \mathbf{x}_i^\top \boldsymbol{\beta}_j, \sigma_j^2) \} K_h(u_i - u) \\ &= \sum_{i=1}^n \sum_{j=1}^C Z_{ij} \left\{ \log \pi_j - \frac{1}{2\sigma_j^2} (y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_j)^2 - \frac{1}{2} \log(\sigma_j^2) \right\} K_h(u_i - u) \end{aligned}$$

Given the current value of the parameter $\boldsymbol{\psi}^{(m)}(u)$, at the $(m+1)$ -th iteration, the EM algorithm proceeds in two steps as follows.

E-step: Since the Z_{ij} 's are unobservable, we compute the expectation of the complete local log-likelihood $\tilde{\ell}_n^c$ with respect to Z_{ij} conditional on the observations $\{(u_i, \mathbf{x}_i, y_i) : i = 1, \dots, n\}$ and the current parameter values $\boldsymbol{\psi}^{(m)}(u)$. This boils down to the computation of the conditional expectations

$$w_{ij}^{(m)} = \frac{\pi_j^{(m)} N(y_i; \mathbf{x}_i^\top \boldsymbol{\beta}_j^{(m)}, \sigma_j^{(m)2})}{\sum_{j=1}^C \pi_j^{(m)} N(y_i; \mathbf{x}_i^\top \boldsymbol{\beta}_j^{(m)}, \sigma_j^{(m)2})}$$

for all $i = 1, \dots, n$ and $j = 1, \dots, C$.

M-step: We maximize the objective function

$$\sum_{i=1}^n \sum_{j=1}^C w_{ij}^{(m)} \left\{ \log \pi_j - \frac{1}{2\sigma_j^2} (y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_j)^2 - \frac{1}{2} \log(\sigma_j^2) \right\} K_h(u_i - u)$$

with respect to $(\pi_j, \boldsymbol{\beta}_j, \sigma_j^2)$. The locally-constant parameter updates are,

$$\begin{aligned}\pi_j^{(m+1)}(u) \equiv \pi_j^{(m+1)} &= \frac{\sum_{i=1}^n w_{ij}^{(m)} K_h(u_i - u)}{\sum_{i=1}^n K_h(u_i - u)} \\ \boldsymbol{\beta}_j^{(m+1)}(u) \equiv \boldsymbol{\beta}_j^{(m+1)} &= (\mathbf{X}^\top \mathbf{W}_j^{(m)} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}_j^{(m)} \mathbf{y} \\ \sigma_j^{2(m+1)}(u) \equiv \sigma_j^{2(m+1)} &= \frac{\sum_{i=1}^n w_{ij}^{(m)} K_h(u_i - u) \left(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_j^{(m+1)} \right)^2}{\sum_{i=1}^n w_{ij}^{(m)} K_h(u_i - u)},\end{aligned}$$

where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$, $\mathbf{y} = (y_1, \dots, y_n)^\top$, and

$$\mathbf{W}_j^{(m)}(u) \equiv \mathbf{W}_j^{(m)} = \text{diag} \left\{ w_{ij}^{(m)} K_h(u_i - u); i = 1, \dots, n \right\}.$$

Web Appendix F.2.2 *The EM algorithm for the MPLLE.* The complete (total) local log-likelihood is given by

$$L_n^c(\boldsymbol{\Psi}; h) = \sum_{t=1}^n \sum_{i=1}^n \sum_{j=1}^C Z_{ij} \left\{ \log \pi_{j,t} + \log f(y_i; \boldsymbol{\theta}_{j,t}(\mathbf{x}_i), \phi_{j,t}) \right\} K_h(u_i - u_t).$$

Given the current value of the parameter $\boldsymbol{\Psi}^{(m)}$, at the $(m+1)$ th iteration, the EM algorithm proceeds in two steps as follows.

E-step: Since the Z_{ij} 's are unobservable, we compute the expectation of the penalized complete local log-likelihood \tilde{L}_n^c with respect to Z_{ij} conditional on the observations $\{(u_i, \mathbf{x}_i, y_i) : i = 1, \dots, n\}$ and the current parameter values $\boldsymbol{\Psi}^{(m)}$. This indeed boils down to the computation of the conditional expectations

$$w_{ij}^{(m)} = \frac{\pi_j^{(m)}(u_i) f(y_i; \boldsymbol{\theta}_j^{(m)}(\mathbf{x}_i, u_i), \phi_j^{(m)}(u_i))}{\sum_{j=1}^C \pi_j^{(m)}(u_i) f(y_i; \boldsymbol{\theta}_j^{(m)}(\mathbf{x}_i, u_i), \phi_j^{(m)}(u_i))}$$

for all $i = 1, \dots, n$ and $j = 1, \dots, C$.

M-step: The objective function \tilde{Q} in M-step of the modified EM algorithm in Section 5 of the paper becomes

$$\begin{aligned}\tilde{Q}(\boldsymbol{\Psi}; \boldsymbol{\Psi}^{(m)}) &= \sum_{j=1}^C \sum_{t=1}^n \left\{ \sum_{i=1}^n w_{ij}^{(m)} \left\{ \log \pi_{j,t} - \frac{1}{2\sigma_{j,t}^2} (y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_{j,t})^2 - \frac{1}{2} \log(\sigma_{j,t}^2) \right\} \right. \\ &\quad \left. \times K_h(u_i - u_t) - \frac{1}{2} \boldsymbol{\beta}_{j,t}^\top \boldsymbol{\Sigma}_j^{(m)} \boldsymbol{\beta}_{j,t} \right\},\end{aligned}$$

where $\Sigma_j^{(m)} = \text{diag}\{\tau_{jl}^{(m)} : l = 1, \dots, d\}$ with

$$\tau_{jl}^{(m)} = \frac{1}{n} \frac{p'_n(\|\mathbf{b}_{jl}^{(m)}\|/\sqrt{n}; \lambda_j)}{\|\mathbf{b}_{jl}^{(m)}\|/\sqrt{n}}.$$

The locally-constant parameter updates are, $t = 1, \dots, n$,

$$\begin{aligned} \pi_{j,t}^{(m+1)} &= \pi_j^{(m+1)}(u_t) = \sum_{i=1}^n w_{ij}^{(m)} K_h(u_i - u_t) / \sum_{i=1}^n K_h(u_i - u_t) \\ \boldsymbol{\beta}_{j,t}^{(m+1)} &= \boldsymbol{\beta}_j^{(m+1)}(u_t) = \arg \min_{\boldsymbol{\beta}_{j,t} \in \mathbb{R}^d} \left\{ \sum_{i=1}^n w_{ij}^{(m)} \left\{ \frac{1}{\sigma_{j,t}^{2(m)}} (y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_{j,t})^2 \right\} K_h(u_i - u_t) \right. \\ &\quad \left. + \boldsymbol{\beta}_{j,t}^\top \Sigma_j^{(m)} \boldsymbol{\beta}_{j,t} \right\} \\ \sigma_{j,t}^{2(m+1)} &\equiv \sigma_j^{2(m+1)}(u_t) = \frac{\sum_{i=1}^n w_{ij}^{(m)} K_h(u_i - u_t) (y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_{j,t}^{(m+1)})^2}{\sum_{i=1}^n w_{ij}^{(m)} K_h(u_i - u)}. \end{aligned}$$

Web Appendix F.3 *The EM algorithm for the MLE and MPLLE in t FM-VCR*

We adapted the EM algorithm outlined in [Yao et al. \(2014\)](#) for t -distribution mixtures to our penalization method for the FM-VCR models.

In a t FM-VCR, the conditional density (mass) of $Y | (\mathbf{X} = \mathbf{x}, U = u)$ is given by

$$f_C^*(y | \boldsymbol{\psi}(u), \mathbf{x}) = \sum_{j=1}^C \pi_j(u) f_T(y - \mathbf{x}^\top \boldsymbol{\beta}_j(u); \sigma_j(u), v_j),$$

where f_T is the density of a t -distribution with v_k degrees of freedom and

$$f_T(\epsilon_j(u); \sigma_j(u), v_j) = \frac{\Gamma\left(\frac{v_j+1}{2}\right)}{(\pi v_j)^{1/2} \Gamma\left(\frac{v_j}{2}\right) \left\{1 + \frac{\epsilon_j^2(u)}{\sigma_j^2(u) v_j}\right\}^{(v_j+1)/2}} \sigma_j(u)$$

with $\epsilon_j(u) = y - \mathbf{x}^\top \boldsymbol{\beta}_j(u)$.

Let $(u_i, \mathbf{x}_i, y_i), i = 1, \dots, n$ be the observations based on a random sample from above model. The (conditional) log-likelihood function is given by

$$\mathcal{L} = \sum_{i=1}^n \log \left\{ \sum_{j=1}^C \pi_j(u_i) f_T(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_j(u_i); \sigma_j(u_i), v_j) \right\}.$$

At any $u \in \mathcal{U}$, the local-kernel log-likelihood in (7) of the paper is given by

$$\ell_n(\boldsymbol{\psi}(u); h) = \sum_{i=1}^n \log \left\{ \sum_{j=1}^C \pi_j f_T(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_j; \sigma_j, v_j) \right\} K_h(u_i - u), \quad (\text{A.75})$$

The locally constant vector of parameters is $\boldsymbol{\psi}(u) = (\boldsymbol{\pi}^\top, \boldsymbol{\phi}^\top, \boldsymbol{\beta}_1^\top, \dots, \boldsymbol{\beta}_C^\top)$, where $\boldsymbol{\pi}^\top =$

$(\pi_1, \dots, \pi_{C-1})$, $\boldsymbol{\phi}^\top = (\sigma_1, \dots, \sigma_C)$ and $\boldsymbol{\beta}_j^\top = (\beta_{j1}, \dots, \beta_{jd})$ for $j = 1, \dots, C$. Note that the entries of the vector $\boldsymbol{\psi}(u)$ are local constant approximations of the functions $\beta_{jl}(u)$, $l = 1, \dots, d$, $\pi_j(u)$, and $\sigma_j(u)$, for all $j = 1, \dots, C$. These entries clearly depend on u , and for simplicity, we suppress u in the notation but keep it for $\boldsymbol{\psi}(u)$.

Using (A.75), the corresponding (total) local-kernel log-likelihood in (8) of the paper is given by

$$\begin{aligned} L_n(\boldsymbol{\Psi}; h) &= \sum_{t=1}^n \ell_n(\boldsymbol{\psi}(u_t); h) \\ &= \sum_{t,i=1}^n \log \left\{ \sum_{j=1}^C \pi_{j,t} f_T(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_{j,t}; \sigma_{j,t}, v_j) \right\} K_h(u_i - u_t), \end{aligned} \quad (\text{A.76})$$

where $\boldsymbol{\Psi}$ is the $(n \times \{C(d+2) - 1\})$ -dimensional matrix of the locally constant vectors $\boldsymbol{\psi}(u_i)$, $i = 1, \dots, n$.

As discussed in Yao et al. (2014), a t -distribution can be represented as a scale mixture of normal distributions as follows. Let g be the latent variable such that

$$\epsilon|g \sim N(0, \sigma^2/g), \quad g \sim \text{Gamma}(v/2, v/2),$$

where $\text{Gamma}(\alpha, \gamma)$ has the density

$$f_G(g; \alpha, \gamma) = \frac{1}{\Gamma(\alpha)} \gamma^\alpha g^{\alpha-1} e^{-\gamma g}, \quad g > 0.$$

Then ϵ has a marginal t -distribution with degrees of freedom v and scale parameter σ . The scale mixture is used in the EM algorithm outlined below.

Web Appendix F.3.1 The EM algorithm. We view the observed data $\{(u_i, \mathbf{x}_i, y_i) : i = 1, \dots, n\}$ as incomplete, and introduce the unobserved Bernoulli random variables Z_{ij} to represent the membership of the i -th observation to the j -th component of the mixture model, $\forall j = 1, \dots, C$, such that $\Pr(Z_{ij} = 1 | u_i, \mathbf{x}_i, y_i) = \pi_j(u_i)$. The complete data consist of $\{(u_i, \mathbf{x}_i, y_i, \mathbf{z}_i, g_i) : i = 1, \dots, n\}$, where $\mathbf{z}_i = (z_{i1}, \dots, z_{iC})^\top$, and g_i 's are the latent variables in the scale mixture.

The complete local-kernel log-Likelihood is given by

$$\begin{aligned}
L_n^c(\Psi; h) &= \sum_{t=1}^n \sum_{i=1}^n \sum_{j=1}^C z_{ij} \log \left\{ \pi_{j,t} N \left(y_i; \mathbf{x}_i^\top \boldsymbol{\beta}_{j,t}, \frac{\sigma_{j,t}^2}{g_i} \right) f_G \left(g_i, \frac{v_j}{2}, \frac{v_j}{2} \right) \right\} K_h(u_i - u_t) \\
&= \sum_{t=1}^n \sum_{i=1}^n \sum_{j=1}^C z_{ij} K_h(u_i - u_t) \log(\pi_{j,t}) \\
&\quad + \sum_{t=1}^n \sum_{i=1}^n \sum_{j=1}^C z_{ij} K_h(u_i - u_t) \log \left\{ N \left(y_i; \mathbf{x}_i^\top \boldsymbol{\beta}_{j,t}, \frac{\sigma_{j,t}^2}{g_i} \right) \right\} \\
&\quad + \sum_{t=1}^n \sum_{i=1}^n \sum_{j=1}^C z_{ij} K_h(u_i - u_t) \log \left\{ f_G \left(g_i, \frac{v_j}{2}, \frac{v_j}{2} \right) \right\}
\end{aligned}$$

The EM algorithm first computes the expected value of $L_n^c(\Psi; h)$, with respect to the latent variables (z_{ij}, g_i) 's conditional on the data, and the current parameter values $\Psi^{(m)}$.

E-step: At the $(m + 1)$ -th iteration, we compute

$$w_{ij}^{(m)} = E(Z_{ij} | \Psi^{(m)}, \text{data}) = \frac{\pi_j^{(m)}(u_i) f_T \left(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_j^{(m)}(u_i); \sigma_j^{(m)}(u_i), v_j \right)}{\sum_{k=1}^C \pi_k^{(m)}(u_i) f_T \left(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_k^{(m)}(u_i); \sigma_k^{(m)}(u_i), v_k \right)}$$

and

$$g_{ij}^{(m)} = E(g_i | \Psi^{(m)}, \text{data}, Z_{ij} = 1) = \frac{v_j + 1}{v_j + \left\{ \frac{y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_j^{(m)}(u_i)}{\sigma_j^{(m)}(u_i)} \right\}^2}.$$

M-step: The MLE parameter updates are given as follows.

$$\pi_{j,t}^{(m+1)} = \pi_j^{(m+1)}(u_t) = \frac{\sum_{i=1}^n w_{ij}^{(m)} K_h(u_i - u_t)}{\sum_{i=1}^n K_h(u_i - u_t)},$$

and

$$\boldsymbol{\beta}_{j,t}^{(m+1)} = \boldsymbol{\beta}_j^{(m+1)}(u_t) = \left(\mathbf{X}^\top W_j^{(m+1)} G_j^{(m+1)} \mathbf{X} \right)^{-1} \mathbf{X}^\top W_j^{(m+1)} G_j^{(m+1)} \mathbf{y},$$

where $W_j^{(m+1)} = \text{diag} \left\{ w_{ij}^{(m+1)}, i = 1, \dots, n \right\}$, $G_j^{(m+1)} = \text{diag} \left\{ g_{ij}^{(m+1)}, i = 1, \dots, n \right\}$, and

$$\sigma_{j,t}^{2(m+1)} = \sigma_j^{2(m+1)}(u_t) = \frac{\sum_{i=1}^n w_{ij}^{(m+1)} g_{ij}^{(m+1)} K_h(u_i - u_t) \left(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_{j,t}^{(m+1)} \right)^2}{\sum_{i=1}^n w_{ij}^{(m+1)} K_h(u_i - u_t)}.$$

The updates of the regression parameters based on ridge penalty are given as

$$\boldsymbol{\beta}_{j,t}^{(m+1)} = \boldsymbol{\beta}_j^{(m+1)}(u_t) = \left(\mathbf{X}^\top W_j^{(m+1)} G_j^{(m+1)} \mathbf{X} + \gamma \mathbf{I}_{d \times d} \right)^{-1} \mathbf{X}^\top W_j^{(m+1)} G_j^{(m+1)} \mathbf{y},$$

where $\gamma > 0$ is the ridge tuning parameter.

Finally, the MPLLE updates of the regression parameters are given as

$$\boldsymbol{\beta}_{j,t}^{(m+1)} = \boldsymbol{\beta}_j^{(m+1)}(u_t) = \left(\mathbf{X}^\top W_j^{(m+1)} G_j^{(m+1)} \mathbf{X} + \gamma \mathbf{I}_{d \times d} + \boldsymbol{\Sigma}_j^{(m+1)} \right)^{-1} \mathbf{X}^\top W_j^{(m+1)} G_j^{(m+1)} \mathbf{y},$$

where $\boldsymbol{\Sigma}_j^{(m)} = \text{diag}\{\tau_{jl}^{(m)} : l = 1, \dots, d\}$ and

$$\tau_{jl}^{(m)} = \frac{1}{n} \frac{p'_n(\|\mathbf{b}_{jl}^{(m)}\|/\sqrt{n}; \lambda_j)}{\|\mathbf{b}_{jl}^{(m)}\|/\sqrt{n}}.$$

Web Appendix F.4 Initial value for the EM algorithm

In our simulation, we obtained $\boldsymbol{\Psi}^{(0)}$ by adding perturbation Gaussian noises to the true value $\boldsymbol{\Psi}^0$. In another approach, we first fit a finite mixture of polynomial regression coefficients of a prespecified degree (five) with fixed (non-varying) mixing probabilities and dispersion parameters (π_j, ϕ_j) . The resulting estimates are then used as the initial value $\boldsymbol{\Psi}^{(0)}$ in the modified EM algorithm. In our simulations, both approaches led to similar results. In our analysis of the real data, we adopted the second approach by the first fitting finite a mixture of polynomial regressions based on several randomly generated initial values.

When the dimension d is comparable to the sample size n , or there is a near-singularity among the covariates, a ridge penalty may be applied in fitting the finite mixture of polynomial regression models.

Web Appendix G. Tuning parameter and mixture order selection

In practice, one needs to choose the band-width h , the tuning parameters λ_j 's, and the mixture order C . It is computationally infeasible to simultaneously choose appropriate values for all these parameters. We now describe data-adaptive strategies for their selection.

Web Appendix G.1 Band-width selection for smoothing

We use the idea of a multi-fold cross-validation (Geisser, 1975) for band-width selection.

Let $\mathcal{D} = \{(u_i, \mathbf{x}_i, y_i), i = 1, \dots, n\}$ represent the full data. We partition \mathcal{D} randomly

into subsets $\mathcal{P}_l, l = 1, \dots, J$, each of size approximately n/J . For each l , we refer to \mathcal{P}_l and $\mathcal{T}_{-l} = \mathcal{D} \setminus \mathcal{P}_l$ as test and training data, respectively. For each l , and given h , let $\tilde{\Psi}_{-l}(h) \equiv \tilde{\Psi}_{-l}$ be the MLLE of Ψ , by maximizing the local-kernel log-likelihood in (7) of the paper and using the training data \mathcal{T}_{-l} . Inspired by Huang et al. (2018), we compute the predictive log-likelihood

$$Pl(h) = \sum_{l=1}^J \sum_{i \in \mathcal{P}_l} \log \left\{ \sum_{j=1}^C \tilde{\pi}_{j,-l}(u_i) f(y_i; \tilde{\theta}_{j,-l}(\mathbf{x}_i, u_i), \tilde{\phi}_{j,-l}(u_i)) \right\}, \quad (\text{A.77})$$

where $\tilde{\theta}_{j,-l}(\mathbf{x}_i, u_i) = g(\mathbf{x}_i^\top \tilde{\beta}_{j,-l}(u_i))$. The predictive log-likelihood is computed over a grid of h values, and we choose a value of h that maximizes $Pl(h)$. To reduce the effect of random partitioning, we repeat the partitioning, say B times, and maximize an average predictive likelihood with respect to h . In our simulation, we used $J = 5$ and $B = 10$.

In what follows, we fix h at the value obtained based on the criterion in (A.77).

Web Appendix G.2 Tuning parameter for variable selection

We use a BIC-type criterion (Wang et al., 2007) for choosing a presumed common tuning parameter $\lambda_1 = \dots = \lambda_C = \lambda$ for the variable selection via penalty p_n in (10) of the paper. Using this method, one may choose different λ_j 's by searching over a C -dimensional grid which is computationally more expensive.

For a value λ , let $\hat{\Psi}_n(\lambda)$ be the MPLLE as in (11) of the paper. Note that to perform variable selection, each non-parametric function $\beta_{jl}(u)$ is estimated at the points u_1, \dots, u_n by the vector $\hat{\mathbf{b}}_{jl}(\lambda) = (\hat{\beta}_{jl}(u_1), \dots, \hat{\beta}_{jl}(u_n))^\top$. Thus, the total number of estimated non-zero regression functions is given by $\sum_{j=1}^C \sum_{l=1}^d 1\{\|\hat{\mathbf{b}}_{jl}(\lambda)\|_2 > 0\}$. We compute

$$\text{BIC}_1(\lambda) = -2\mathcal{L}(\hat{\Psi}_n(\lambda)) + \log n \times \left\{ \sum_{j=1}^C \sum_{l=1}^d 1\{\|\hat{\mathbf{b}}_{jl}(\lambda)\|_2 > 0\} \right\} \times \text{DF}_h, \quad (\text{A.78})$$

where \mathcal{L} is the log-likelihood in (6) of the paper evaluated at $\hat{\Psi}_n(\lambda)$, and as in Huang et al. (2013),

$$\text{DF}_h = \tau_K h^{-1} |\mathcal{U}| \left\{ K(0) - \frac{1}{2} \int_{\mathcal{U}} K^2(u) du \right\},$$

where $|\mathcal{U}|$ is equal to the length of the support of the index variable U , and

$$\tau_K = \left\{ K(0) - \frac{1}{2} \int_{\mathcal{U}} K^2(u) du \right\} \left\{ \int_{\mathcal{U}} \left\{ K(u) - \frac{1}{2} K * K(u) \right\}^2 du \right\}^{-1}$$

and $K * K(t) = \int_0^t K(u)K(t-u)du$. In [Web Appendix G.4](#) that follows, we compute the degrees of freedom DF_h for the Epanechnikov Kernel. We compute the BIC over a grid of λ -values, say $\lambda_1, \dots, \lambda_M$, and then choose a value of λ that minimizes [\(A.78\)](#).

In our simulation and data analysis, guided by the theory in [Section 4](#) of the paper, for a given n , the smoothing parameter h is chosen by maximizing the predictive log-likelihood [\(A.77\)](#) over the range $[0.1, 2n^{-0.2}]$; the tuning parameter λ is chosen by minimizing the BIC in [\(A.78\)](#) over the range $[10^{-6}, 3n^{-0.45}]$ for the AdpLASSO and over the range $[10^{-6}, 3n^{-0.35}]$ for the LASSO, SCAD, MCP. The constants involved in the ranges are chosen by trial and error.

Web Appendix G.3 Mixture order selection

We also use a BIC for selection of the mixture order C , when it is unknown. Here, we need to take into account the total number of non-parametric mixing probabilities and dispersion parameters that are estimated. We compute the BIC

$$\text{BIC}_2(C) = -2\mathcal{L}(\widehat{\Psi}_n(C)) + \log n \times \left\{ 2C - 1 + \sum_{j=1}^C \sum_{l=1}^d 1_{\{\|\widehat{\mathbf{b}}_{jl}(C)\|_2 > 0\}} \right\} \times \text{DF}_h, \quad (\text{A.79})$$

where $\widehat{\Psi}_n(C)$ is the MPLLE obtained for the fitted mixture models of different orders $C = 1, \dots, \mathcal{K}$, for some pre-specified upper bound \mathcal{K} . In [Web Appendix G.4](#) that follows, we compute the degrees of freedom DF_h for the Epanechnikov Kernel. The mixture order is then estimated by \widehat{C}_n that minimizes $\text{BIC}_2(C)$ over $1 \leq C \leq \mathcal{K}$. Theoretically, under the conditions of [Theorem 2](#) in the paper, and similar to the results of [Leroux \(1992\)](#), the estimator \widehat{C}_n does not underestimate the true mixture order with probability tending to one as $n \rightarrow \infty$. Our simulation results show that \widehat{C}_n does not underestimate the mixture order, and also the percentage of overestimation decreases as n increases.

Web Appendix G.4 *Computing the degree of freedom for the BIC*

For the Epanechnikov Kernel in [Web Appendix C](#), we have

$$\text{DF}_h = \tau_K h^{-1} |\mathcal{U}| \left\{ K(0) - \frac{1}{2} \int_{\mathcal{U}} K^2(u) du \right\},$$

and for the BIC given in [Web Appendix G](#) of the paper, where \mathcal{U} is the support of the index variable U , $|\mathcal{U}|$ is equal to the length of the support, and

$$\tau_K = \frac{K(0) - \frac{1}{2} \int K^2(u) du}{\int_{\mathcal{U}} \left\{ K(u) - \frac{1}{2} K * K(u) \right\}^2 du},$$

with the convolution

$$K * K(t) = \int_0^t K(u) K(t-u) du.$$

For the Epanechnikov Kernel, we have

$$K * K(t) = \frac{3^2}{4^2} \int_0^t \frac{3}{4} (1-u^2) \frac{3}{4} (1-(t-u)^2) du = \frac{3^2}{4^2} \left(t - \frac{2}{3} t^3 + \frac{1}{30} t^5 \right),$$

and

$$\begin{aligned} \int_0^1 \left\{ K(u) - \frac{1}{2} K * K(u) \right\}^2 dt &= \int_0^1 \left\{ \frac{3}{4} (1-u^2) - \frac{1}{2} \frac{3^2}{4^2} \left(u - \frac{2}{3} u^3 + \frac{1}{30} u^5 \right) \right\}^2 du \\ &= \frac{3^2}{4^2} \int_0^1 \left(1 - \frac{3}{8} u - u^2 + \frac{1}{4} u^3 - \frac{1}{80} u^5 \right)^2 du \\ &= 0.2285. \end{aligned}$$

Thus,

$$K(0) = \frac{3}{4},$$

$$\int_0^1 K^2(u) du = \int_0^1 \frac{3^2}{4^2} (1-u^2)^2 du = \frac{3}{10},$$

and

$$\text{DF}_h = \frac{3/4 - 3/20}{0.2285} \frac{|\mathcal{U}|}{h} (3/4 - 3/20) = \frac{1.5755}{h}.$$

Web Appendix H. Additional simulation studies

Web Appendix H.1 *Model Misspecification*

Since our proposed estimation and variable selection method is likelihood-based, one expects that a misspecification of the parametric form of the mixture components would, in general, affect the performance of the method, although the degree of performance degradation depends on how different the misspecified and the true models are. In terms of the number of latent classes, if C is smaller than its true value, then the fitted model converges to an under-specified model that minimizes the Kullback-Leibler distance to the true FM-VCR model (Leroux, 1992). Thus, we cannot make statements about the consistent estimation of the parameters of the true model based on an under-specified model. On the other hand, if C is larger than its true value, then the overall behavior of the true FM-VCR model will be captured by the over-fitted mixture model. In particular, Ho and Nguyen (2016) and Ho et al. (2022) showed that the density function of an over-fitted standard finite mixture model and the finite mixture of regressions (both special cases of FM-VCRs) consistently estimates the true finite mixture density. We believe such property also holds for over-fitted FM-VCR, although further investigation is needed to study this theoretically.

Below we investigate, through simulation, two scenarios related to model misspecification: misspecifying the parametric form and misspecifying the number of components, with the results presented below as (i) and (ii), respectively. Our simulations are based on the dimension $d = 10$ and sample sizes $n = 200$ and 400 , representative settings of those in the paper. The results are based on $R = 100$ replicates.

(i) Misspecification of the parametric form of the mixture component density

We generated random samples $(u_i, \mathbf{x}_i, y_i), i = 1, \dots, n$, from a two-components FM-VCR model with each component density set to be a t -distribution with 10 degrees of freedom,

which has moderately heavier tails than a Gaussian distribution. The rest of the parameter settings are the same as those in Table 1 of the Simulation Section in the paper.

Using our proposed regularization method outlined in Section 3 of the paper, we fitted two models to the generated data: the correct t FM-VCR from which the data were generated and the misspecified Gaussian FM-VCR model. As in the paper, the simulation results are summarized in terms of sensitivity, specificity, and estimation error (L_2). The results are given in Web Table 7 below. Similar to the Gaussian FM-VCR model, the details of the numerical implementation of the penalization method for a t FM-VCR are now given in [Web Appendix F.3](#) above.

From Web Table 7, we observe that, overall, the results based on the correct t FM-VCR (upper portion of the table) are reasonable in terms of all the performance measures under consideration. The results are roughly similar to those of Table 2 of the paper based on the correct Gaussian FM-VCR model. As expected, the performance of the method in Component 2 (the larger component) of the mixture is better. On the other hand, the results (lower portion of the table) based on the misspecified Gaussian model show that the performance of the method in terms of sensitivity and specificity generally degrades but not dramatically in terms of sensitivity, while specificity has not been affected much. The quality of the parameter estimates in terms of estimation error has been affected particularly in Component 1, the smaller component. As the sample size n increases, the loss in performance is less which is expected. Overall, the misspecification for the above parameter setting has not affected the performance much.

(ii) Misspecification of the number of mixture components (C)

We consider the Gaussian FM-VCR with the true number of components (order) $C = 2$ and the same parameter setting as in Table 1 of the paper. We generated random samples $(u_i, \mathbf{x}_i, y_i), i = 1, \dots, n$, from this model and fitted Gaussian FM-VCR models with orders

$C = 1, 2, 3$ and 4 , to the data. For the overfitted FM-VCR models with orders $C = 3, 4$, we first perform a component-matching in which we find the closest components of the overfitted model to those of the true model with $C = 2$ and then compute the performance measures.

The simulation results are summarized in Web Table 8. The results for the fitted misspecified models with orders $C = 1, 3$ and 4 are compared with those from the correct order $C = 2$. The one-component underfitted model with $C = 1$ resembles the behavior of the larger component of the correct model, i.e. Component 2, but with lower sensitivity and specificity and larger estimation errors for the corresponding measures for Component 2 when $C = 2$ was fitted. On the other hand, the behavior of the models with $C = 3, 4$ are similar to those with the correct model but with lower estimation errors. These simulation results are in line with the theoretical properties of over-fitted finite mixture models (Ho and Nguyen, 2016; Ho et al., 2022), although further work is necessary to theoretically investigate such properties for the over-fitted FM-VCR models.

Web Appendix H.2 Comparison of FMR and FM-VCR models

We generated data from a Gaussian finite mixture of regression (FMR) model (without allowing for varying coefficients) with the number of components $C = 2$, dimension $d = 10$, and samples sizes $n = 200$ or 400 . The parameter setting for the FMR model is as follows:

$$\beta_1^\top = [-0.5, 0.25, 0.25, 0, -0.25, 0, 0, 0, 0, 0, 0],$$

$$\beta_2^\top = [-0.25, 0.25, 0, 0.25, 0, 0, 0, 0, 0, 0, 0],$$

$$\sigma^\top = [0.39, 0.45], \quad \pi^\top = [0.55, 0.45].$$

We then fitted both Gaussian FMR and FM-VCR models to each simulated sample. The results are given in Web Table 9. We observe that, for sensitivity, FM-VCR (the wrong model) is generally worse (on average based on 100 replicates) than FMR (the true underlying

model), with larger standard deviations (SD). On the other hand, for specificity, FM-VCR is consistently better and with much smaller SD. This implies that FM-VCR selected fewer variables, hence fewer true positives and fewer false positives, thus larger specificity and smaller sensitivity. This result makes sense because with a more complex model and more parameters but a fixed sample size, variables are harder to be selected as significant, especially for the parameter setting considered above where the effects of covariates (β_{jl} 's) are weak. Thus, in practice, one may fit both models to a dataset and assess the results as we have done for our real data.

Web Appendix H.3 *Alternative tuning parameter selection*

In general, it is challenging to provide a universally dominant method for a data-dependent choice of the three tuning parameters. In addition to the implementation of the sequential selection of h and then λ for a prespecified mixture order C (Web Appendix G), here we investigate an alternative tuning parameter selection approach, where h and λ are selected simultaneously, to see whether performance can be improved although at the expense of greater computational intensity. Specifically, we investigate the scenario in which $C = 2$, dimension $d = 10$, and sample sizes $n = 200, 400$, which was considered in the paper. The variable selection and estimation error results are given in Web Table 10, and the time (in seconds) are given in Web Table 11.

As we can see from the computational time results, on average, the MPLLE based on selecting h and then λ is at least 18 times faster than the MPLLE when (h, λ) are selected simultaneously, which is expected. On the other hand, comparing the results of Web Table 10 with those in Table 2 of the paper, despite the increase in computational time, the gain in performance of the MPLLE based on the simultaneous selection of (h, λ) both in terms of variable selection and estimation error is negligible.

Web Appendix H.4 *High-dimensions*

In this section, we assess the performance of the penalization method when the dimension d exceeds the sample size n , although this case is not covered in our theoretical results where d is fixed as n grows.

We have considered the two-component ($C = 2$) Gaussian FM-VCR model with the parameter setting given in Table 1 of the paper, dimension $d = 500$, and sample sizes $n = 200, 400$. In this model, covariates (x_1, x_2, x_4) in Component 1 and covariates (x_1, x_3) in Component 2 have non-zero $\beta_{jk}(u)$, and the rest have zero coefficients. There are a total of $p = Cd + 2C - 1 = 1003$ non-parametric functions to be estimated by the proposed penalization method. The results are summarized in Web Table 12.

The results show that while the specificity remains high and similar to dimensions $d = 5, 10, 20$, and 50 already considered in Example 1 of Section 6 in the main paper (Table 2, and see also Web Table 2), there is some reduction in sensitivity, although the difference is quite small when compared to $d = 50$, indicating that reasonable results are likely to be obtained by the method even when the number of variables exceed the sample size.

Nevertheless, rigorous study of high-dimensional setting requires new theoretical and numerical tools, which are not covered in the current paper.

In the second part of our simulation to investigate high-dimensional settings, we examined two popular screening techniques that one could use when fitting an FM-VCR to a dataset.

(1) Screening using the Pearson’s Correlation.

We computed the sample correlation between the response and $d = 500$ covariates for each simulated sample, and we kept $d^* = 50$ top covariates with the highest correlation values. This procedure was repeated for 100 replicated samples. Web Table 13 shows the top 20 covariates that have the highest frequencies of being selected out of the 100 replications. Note that in our simulation, the covariates (x_1, x_2, x_3, x_4) have non-zero $\beta_{jl}(u)$ ’s in the true

data generating FM-VCR, and they all survived the screening process in all 100 replicates. Some other covariates (those with zero $\beta_{jl}(u)$'s, essentially noise) also survived with high frequencies, necessitating the use of our FM-VCR method for further in-depth variable selection and estimation.

(2) Screening using marginal likelihood of single-covariate FM-VCR model with $C = 2$.

In this approach, via MLE, we fitted $d = 500$ single-covariate FM-VCR models each of which includes intercepts and only one covariate in both components of the mixture model. We then rank the top $d^* = 50$ covariates based on the likelihood value of their corresponding fitted model. This procedure was repeated for 100 replicated samples. Web Table 14 below shows the top 20 covariates with the highest frequencies of being among the top 50's out of the 100 replications. Again, we can see that the covariates (x_1, x_2, x_3, x_4) that have non-zero $\beta_{jl}(u)$'s in the true data generating FM-VCR survived the screening process in all 100 replicates, among other (noise) covariates that were in fact selected less.

Web Appendix I. More on the data analysis from the OCN study

To substantiate the need for an FM-VCR model, we carried out an analysis by fitting an FMR model to the OCN data to show the lack of fit of this model compared to the FM-VCR model. Specifically, using all four penalties studied in the paper, we analyzed the OCN data by fitting FMR models with mixture orders $C = 1, 2, 3, 4$, and 5. An FMR model with two components ($C = 2$), and based on the SCAD penalty, was selected by the BIC. Similar to the fitted sparse FM-VCR model, SNP rs7456421 was also selected in the fitted sparse FMR model, among several other covariates. The mixing proportion of Component 1 in the FMR model is 53% which is approximately equal to the average of the varying mixing proportion of Component 1 in the FM-VCR model. The variance of Component 1 in the FMR model is about 62, which is close to the average of varying variance of Component 1 in the FM-VCR

model. However, the second component of the FMR model has a much larger variance, about 149, compared to the variance of Component 2 in the FM-VCR model (which ranges from 13 to 19 as a function of age). This much larger variance of Component 2 compared to the variance of Component 1 in the FMR can be an indication of a misspecified model. To demonstrate the lack of fit using the FMR model and the subsequent improvement by the FM-VCR model, we analyzed their residuals by plotting the quantiles against their normal counterparts. As can be seen from the QQ-plots in Web Figure 2 below, the use of FM-VCR indeed improves the fit to the data compared to FMR, especially for Component 2.

We have also supplied two plots to provide further evidence of the varying coefficient effect of SNP rs7456421 over the age on osteocalcin. In Web Figure 3(a), we have added the point-wise error bars to the estimated varying coefficient plot of the SNP, where the error bars were calculated using the EM algorithm approximation method of ?. This substantiates our conclusion that SNP rs7456421 has an age-dependent effect on osteocalcin, and that this SNP has a significantly negative effect at a younger age, and this effect diminishes as one ages. Web Figure 3(b) is the derivative curve (approximated using Matlab) of the estimated coefficient curve over age. It is evident from the plot that the underlying true derivatives are unlikely to be constant at 0 over the age range.

Web Appendix J. Additional tables and figures

[Table 1 about here.]

[Table 2 about here.]

[Table 3 about here.]

[Table 4 about here.]

[Table 5 about here.]

[Table 6 about here.]

[Table 7 about here.]

[Table 8 about here.]

[Table 9 about here.]

[Table 10 about here.]

[Table 11 about here.]

[Table 12 about here.]

[Table 13 about here.]

[Table 14 about here.]

[Figure 1 about here.]

[Figure 2 about here.]

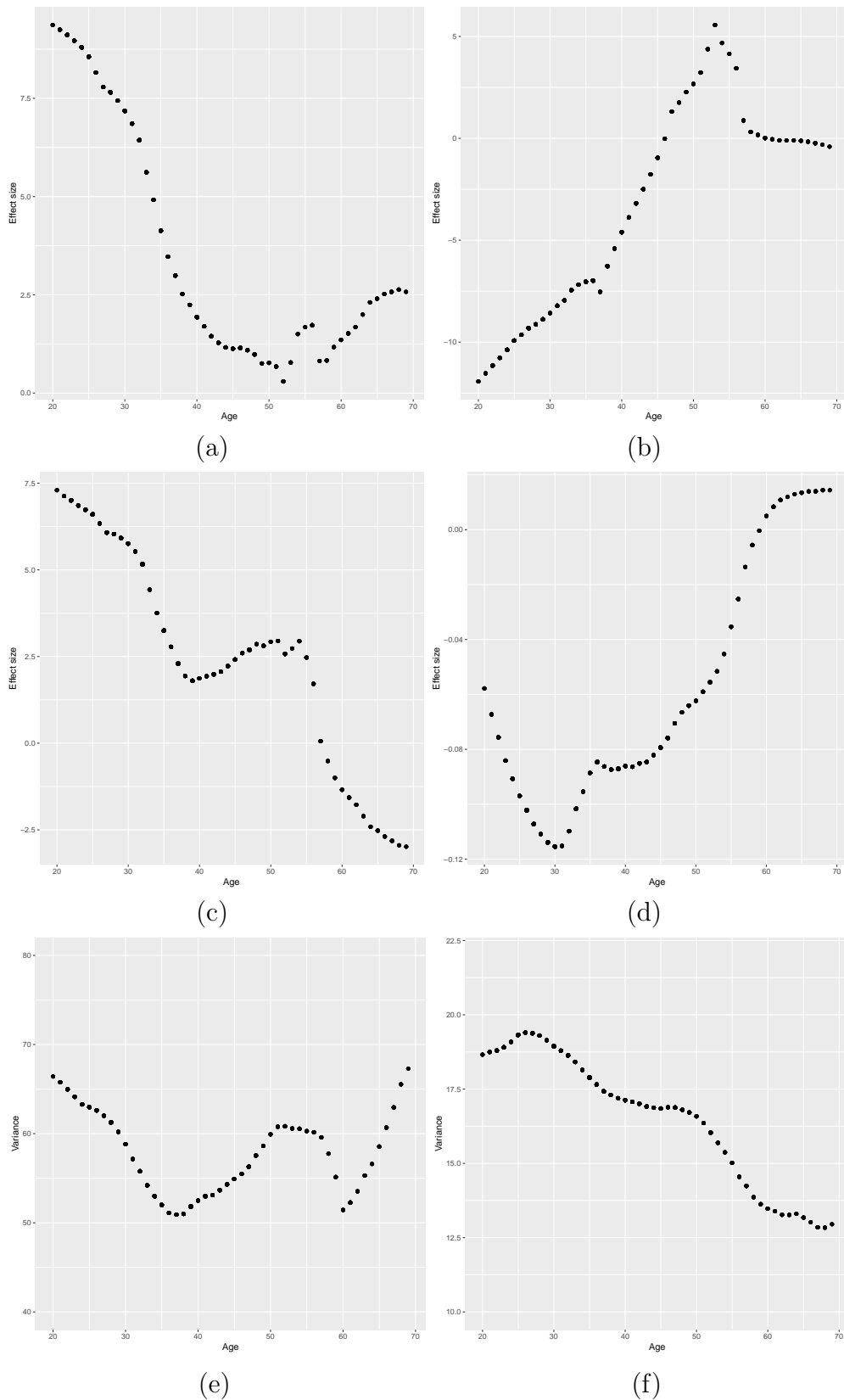
[Figure 3 about here.]

References

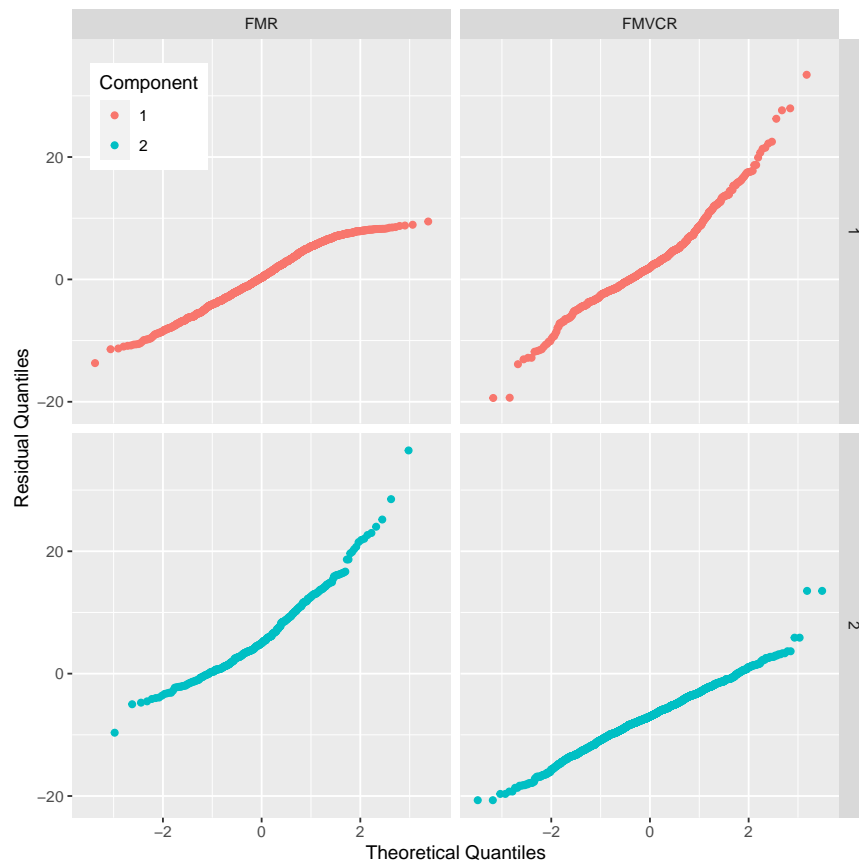
- Berk, R., Brown, L., Buja, A., Zhang, K., and Zhao, L. (2013). Valid post-selection inference. *The Annals of Statistics* **41**, 802 – 837.
- Fang, F., Jiwei, Z., S. Ejaz, A., and Qu, A. (2021). A weak-signal-assisted procedure for variable selection and statistical inference with an informative subsample. *Biometrics* **77**, 996–1010.
- Geisser, S. (1975). The predictive sample reuse method with applications. *Journal of the American Statistical Association* **70**, 320–328.
- Ho, N. and Nguyen, X. (2016). On strong identifiability and convergence rates of parameter estimation in finite mixtures. *Electronic Journal of Statistics* **10**, 271–307.
- Ho, N., Yang, C.-Y., and Jordan, M. I. (2022). Convergence rates for gaussian mixtures of experts. *Journal of Machine Learning Research* .

- Huang, M., Li, R., and Wang, S. (2013). Nonparametric mixture of regression models. *Journal of the American Statistical Association* **108**, 929–941.
- Huang, M., Yao, W., Wang, S., and Chen, Y. (2018). Statistical inference and applications of mixture of varying coefficient models. *Scandinavian Journal of Statistics* **45**, 618–643.
- Janson, S. (1987). Maximal spacing in several dimensions. *Annals of Probability* **15**, 274–80.
- Leroux, B. (1992). Consistent estimation of a mixing distribution. *The Annals of Statistics* **20**, pp. 1350–1360.
- van der Vaart, A. (1998). *Asymptotic Statistics*. Cambridge University Press, Cambridge.
- van der Vaart, A. and Wellner, J. (1996). *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Series in Statistics. Springer.
- Wang, H., Li, R., and Tsai, C.-L. (2007). Tuning parameter selectors for the smoothly clipped absolute deviation method. *Biometrika* **94**, 553–568.
- Wei, F., Huang, J., and Li, H. (2011). Variable selection and estimation in high-dimensional varying-coefficient models. *Statistica Sinica* **21**, 1515–1540.
- Yao, W., Wei, Y., and Yu, C. (2014). Robust mixture regression using the t-distribution. *Computational Statistics & Data Analysis* **71**, 116–127.
- Zhang, D., Khalili, A., and Asgharian, M. (2022). Post-model-selection inference in linear regression models: An integrated review. *Statistics Surveys* **16**, 86–136.
- Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* **101**, 1418–1429.

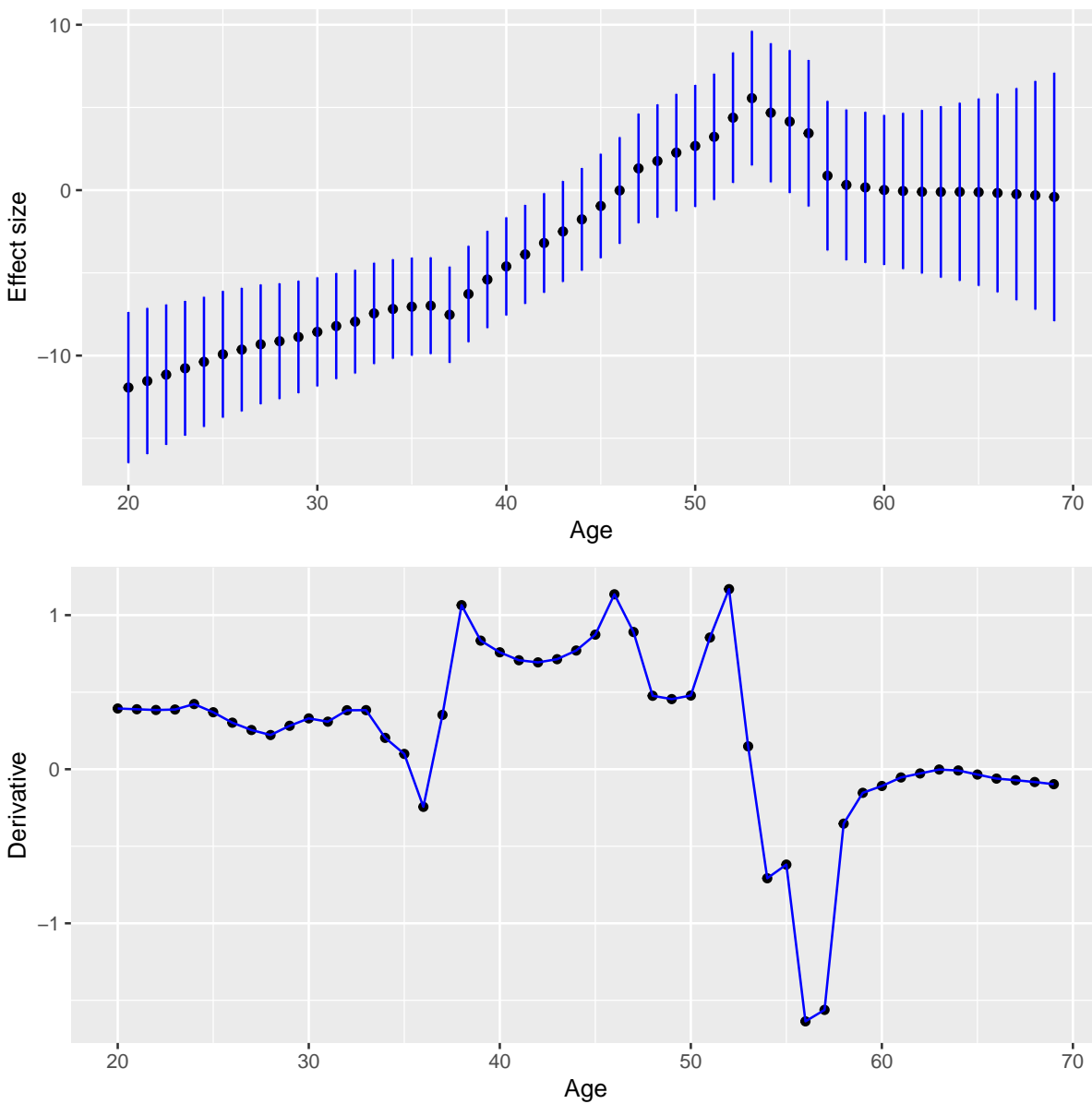
Received January 2022. Revised ? 2022. Accepted ? 2022.



Web Figure 1: Osteocalcin data analysis; (a) The estimated effect of rs109522346 over time in hOCN. (b) The estimated effect of rs7456421 over time in hOCN. (c) The estimated effect of rs4074826 over time in hOCN. (d) The estimated effect of folic acid over time in hOCN. (e) Estimated variance over time in hOCN. (f) Estimated variance over time in lOCN.



Web Figure 2: QQ-plots of the OCN data analyses; results of FM-VCR versus FMR model.



Web Figure 3: The non-constant effect of SNP rs7456421 on osterocalcin for Component 1 of the two-component FM-VCR model. (a) Error bar of the estimated effect size at each age; (b) Derivatives of the non-constant effect curve.

Web Table 1: Parameters settings for the Gaussian FM-VCR models with $C = 2, 3$.

Component(j)	1	2	3
Parameters	$d(= 5, 10, 20, 50)^1$		
$\beta_{j0}(u)$	-2	-1	1
$\beta_{j1}(u)$	$1 + 0.5 \cos(\pi u)$	$1.5 \sin(\pi u)$	0
$\beta_{j2}(u)$	$1 + 0.5 \cos(2\pi u)$	0	$-1.5(\cos(2\pi u))^2$
$\beta_{j3}(u)$	0	$1.5 - 0.5 \sin(\pi u/2)$	0
$\beta_{j4}(u)$	$\sin(6\pi u)$	0	$1 - 0.5 \cos(3u\pi/2)$
$\beta_{j5}(u)$	0	0	0
\vdots	\vdots	\vdots	\vdots
$\beta_{j,50}(u)$	0	0	0
$\sigma_j(u)$	$0.3e^{(0.5u)}$	$0.5e^{(-0.2u)}$	$0.8e^{(0.5u)}$
$\pi_j(u)$	$e^{0.5u}/(1 + e^{0.5u})$	$(1 + e^{0.5u})^{-1}$	—
$\pi_j(u)$	$e^{0.5u}/(1 + e^{0.5u} + e^{0.2u})$	$e^{0.2u}/(1 + e^{0.5u} + e^{0.2u})$	$(1 + e^{0.5u} + e^{0.2u})^{-1}$

¹ We keep the non-zero coefficient functions $\beta_{jl}(\cdot)$ the same for different values of the dimension d .

Web Table 2: Results of Example 1: average (SD) sensitivity, specificity, and estimation errors: $C = 2$.

Criteria	Sensitivity		Specificity		$L_2(\hat{\beta}_j)$		$L_2(\hat{\sigma}_j^2)$		$L_2(\hat{\pi}_j)$	
	1	2	1	2	1	2	1	2	1	2
d = 5 n = 200										
Oracle	—	—	—	—	0.6504	0.1619	0.3828	0.1037	0.3209	0.1037
MLE	—	—	—	—	1.0348	0.2431	0.7476	0.1273	0.3702	0.1273
AdpLASSO	0.6867 _(0.2030)	0.9483 _(0.1914)	0.9667 _(0.1249)	0.9622 _(0.1221)	0.9337	0.2077	0.6119	0.1182	0.5001	0.1182
LASSO	0.7011 _(0.2230)	0.8467 _(0.3467)	0.8133 _(0.2653)	0.8244 _(0.2136)	0.9829	0.2560	0.8119	0.1670	0.5237	0.1670
MCP	0.6822 _(0.1269)	0.9517 _(0.1535)	0.9783 _(0.1099)	0.9722 _(0.1203)	0.8885	0.1632	0.5608	0.0938	0.5046	0.0938
SCAD	0.6844 _(0.1266)	0.9567 _(0.1467)	0.9800 _(0.1063)	0.9744 _(0.1176)	0.8819	0.1655	0.5536	0.0926	0.5074	0.0926
n = 400										
Oracle	—	—	—	—	0.3987	0.0831	0.2144	0.0800	0.2730	0.0800
MLE	—	—	—	—	0.5467	0.0969	0.3895	0.0749	0.2782	0.0749
AdpLASSO	0.9778 _(0.0833)	0.9967 _(0.0577)	0.9983 _(0.0289)	0.9967 _(0.0430)	0.4991	0.0819	0.2253	0.0743	0.3114	0.0743
LASSO	0.9889 _(0.0599)	0.9950 _(0.0645)	0.9383 _(0.1793)	0.9100 _(0.1799)	0.5819	0.1165	0.3759	0.0859	0.3377	0.0859
MCP	0.9833 _(0.0728)	1.0000 _(0.0000)	1.0000 _(0.0000)	0.9989 _(0.0192)	0.4582	0.0777	0.2079	0.0735	0.3128	0.0735
SCAD	0.9844 _(0.0704)	1.0000 _(0.0000)	1.0000 _(0.0000)	0.9989 _(0.0192)	0.4581	0.0774	0.2074	0.0734	0.3129	0.0734
d = 10 n = 200										
Oracle	—	—	—	—	0.6983	0.1616	0.3968	0.0887	0.3143	0.0887
MLE	—	—	—	—	1.2296	0.2886	1.0160	0.1378	0.3355	0.1378
AdpLASSO	0.5889 _(0.2390)	0.8933 _(0.2592)	0.9705 _(0.0667)	0.9663 _(0.0610)	1.0894	0.2876	0.8483	0.1608	0.5007	0.1608
LASSO	0.4700 _(0.3206)	0.7750 _(0.4009)	0.9276 _(0.1002)	0.9167 _(0.0869)	1.3206	0.3829	1.1522	0.2940	0.5054	0.2940
MCP	0.5811 _(0.1940)	0.8950 _(0.2450)	0.9838 _(0.0454)	0.9812 _(0.0481)	1.0420	0.2378	0.7633	0.1056	0.5034	0.1056
SCAD	0.5944 _(0.1878)	0.8917 _(0.2435)	0.9833 _(0.0488)	0.9800 _(0.0503)	1.0294	0.2372	0.7612	0.1066	0.5010	0.1066
n = 400										
Oracle	—	—	—	—	0.4393	0.0888	0.2270	0.0816	0.2728	0.0816
MLE	—	—	—	—	0.8776	0.1489	0.6068	0.0867	0.3374	0.0867
AdpLASSO	0.6900 _(0.1272)	0.9933 _(0.0574)	0.9971 _(0.0232)	0.9954 _(0.0294)	0.8079	0.1654	0.3905	0.0692	0.4110	0.0692
LASSO	0.7656 _(0.2302)	0.9700 _(0.1659)	0.9705 _(0.0657)	0.9617 _(0.0638)	0.8617	0.2322	0.5176	0.1076	0.4249	0.1076
MCP	0.6922 _(0.1042)	0.9917 _(0.0641)	0.9971 _(0.0200)	0.9983 _(0.0144)	0.7892	0.1455	0.3745	0.0616	0.4119	0.0616
SCAD	0.6878 _(0.0979)	0.9917 _(0.0641)	0.9962 _(0.0231)	0.9979 _(0.0160)	0.7919	0.1463	0.3768	0.0620	0.4122	0.0620
d = 20 n = 200										
Oracle	—	—	—	—	0.6147	0.0952	0.2975	0.0790	0.3133	0.0790
MLE	—	—	—	—	1.3671	0.2255	1.3035	0.0821	0.3835	0.0821
AdpLASSO	0.6456 _(0.1153)	0.9433 _(0.2025)	0.9814 _(0.0434)	0.9839 _(0.0416)	0.9140	0.2692	0.7131	0.0563	0.4878	0.0563
LASSO	0.6567 _(0.1200)	0.9533 _(0.1870)	0.9629 _(0.0587)	0.9524 _(0.0532)	0.9240	0.3363	0.8963	0.0563	0.4940	0.0563
MCP	0.6433 _(0.1212)	0.9533 _(0.1986)	0.9941 _(0.0278)	0.9955 _(0.0218)	0.9149	0.2200	0.6193	0.0520	0.4673	0.0520
SCAD	0.6411 _(0.1268)	0.9467 _(0.2061)	0.9943 _(0.0277)	0.9957 _(0.0211)	0.9183	0.2270	0.6333	0.0524	0.4721	0.0524
n = 400										
Oracle	—	—	—	—	0.3533	0.0721	0.2108	0.0670	0.2515	0.0670
MLE	—	—	—	—	0.9120	0.1019	0.6161	0.0377	0.4234	0.0377
AdpLASSO	0.6678 _(0.0192)	0.9978 _(0.0385)	0.9986 _(0.0148)	0.9996 _(0.0048)	0.7984	0.1697	0.4171	0.0356	0.4844	0.0356
LASSO	0.7156 _(0.1181)	1.0000 _(0.0000)	0.9780 _(0.0351)	0.9694 _(0.0412)	0.8226	0.2373	0.5234	0.0365	0.4995	0.0365
MCP	0.6667 _(0.0472)	0.9967 _(0.0577)	0.9994 _(0.0102)	0.9996 _(0.0068)	0.7962	0.1589	0.4331	0.0360	0.4808	0.0360
SCAD	0.6678 _(0.0510)	0.9967 _(0.0577)	0.9992 _(0.0107)	0.9996 _(0.0068)	0.7961	0.1588	0.4331	0.0360	0.4806	0.0360
d = 50 n = 200										
Oracle	—	—	—	—	0.7270	0.1672	0.3934	0.0968	0.2946	0.0968
MLE	—	—	—	—	2.0960	0.2002	2.0600	0.0808	0.1598	0.0808
AdpLASSO	0.5478 _(0.2188)	0.7633 _(0.3307)	0.9882 _(0.0150)	0.9895 _(0.0157)	1.2587	0.2945	1.0778	0.1649	0.4606	0.1649
LASSO	0.5067 _(0.2852)	0.6983 _(0.4148)	0.9873 _(0.0147)	0.9872 _(0.0154)	1.3616	0.3737	1.2348	0.2622	0.4830	0.2622
MCP	0.5567 _(0.1617)	0.8217 _(0.2468)	0.9923 _(0.0105)	0.9925 _(0.0107)	1.2601	0.2437	1.0173	0.1137	0.3995	0.1137
SCAD	0.5578 _(0.1589)	0.8283 _(0.2413)	0.9930 _(0.0102)	0.9926 _(0.0109)	1.2481	0.2397	1.0028	0.1117	0.4045	0.1117
n = 400										
Oracle	—	—	—	—	0.4701	0.1410	0.3067	0.0957	0.2988	0.0957
MLE	—	—	—	—	1.8860	0.2506	1.6380	0.1881	0.2752	0.1881
AdpLASSO	0.6122 _(0.1669)	0.8783 _(0.2299)	0.9906 _(0.0152)	0.9926 _(0.0125)	1.1289	0.2968	0.8221	0.0953	0.5259	0.0953
LASSO	0.6067 _(0.2234)	0.8750 _(0.2590)	0.9887 _(0.0156)	0.9891 _(0.0164)	1.1621	0.3373	0.8745	0.1454	0.5231	0.1454
MCP	0.6067 _(0.1593)	0.8900 _(0.2075)	0.9950 _(0.0092)	0.9953 _(0.0089)	1.0893	0.2396	0.7811	0.0676	0.5090	0.0676
SCAD	0.5978 _(0.1556)	0.8833 _(0.2118)	0.9948 _(0.0094)	0.9951 _(0.0090)	1.1050	0.2428	0.7983	0.0688	0.5053	0.0688

Web Table 5: Results of Example 2: Average (SD) sensitivity, specificity, and estimation errors: ($C = 2, d = 5, n = 200$).

Criteria	Sensitivity		Specificity		$L_2(\hat{\beta}_j)$		$L_2(\hat{\sigma}_j^2)$		$L_2(\hat{\pi}_j)$	
	1	2	1	2	1	2	1	2	1	2
p = 0.05										
Oracle	—	—	—	—	1.1871	0.0771	0.7560	0.0919	0.0765	0.0919
MLE	—	—	—	—	1.5363	0.0678	1.3980	0.0926	0.0740	0.0926
AdpLASSO	0.0878 _(0.1831)	0.5883 _(0.3630)	0.8933 _(0.2092)	0.9889 _(0.0763)	1.7500	0.1263	1.2046	0.1338	0.1192	0.1338
LASSO	0.0822 _(0.1963)	0.6117 _(0.4090)	0.8983 _(0.2057)	0.9578 _(0.1324)	1.7360	0.1348	1.2040	0.1390	0.1189	0.1390
MCP	0.4211 _(0.3000)	0.6517 _(0.2969)	0.8917 _(0.2103)	0.9467 _(0.1340)	1.5986	0.0888	1.2499	0.1029	0.0956	0.1029
SCAD	0.4311 _(0.2975)	0.6600 _(0.2995)	0.8933 _(0.2092)	0.9456 _(0.1264)	1.5994	0.0903	1.2441	0.1038	0.0950	0.1038
p = 0.1										
Oracle	—	—	—	—	0.9970	0.0825	0.6118	0.0724	0.0991	0.0724
MLE	—	—	—	—	1.3279	0.0844	1.0951	0.0863	0.0976	0.0863
AdpLASSO	0.2611 _(0.2917)	0.7367 _(0.3307)	0.8950 _(0.2198)	0.9667 _(0.1377)	1.6326	0.1669	1.0047	0.1540	0.1623	0.1540
LASSO	0.3211 _(0.3375)	0.6983 _(0.3877)	0.8667 _(0.2532)	0.8889 _(0.2152)	1.6067	0.1777	1.1127	0.1582	0.1578	0.1582
MCP	0.5511 _(0.2833)	0.8167 _(0.2613)	0.9067 _(0.1951)	0.9667 _(0.1107)	1.3219	0.0977	0.9244	0.0872	0.1384	0.0872
SCAD	0.5544 _(0.2853)	0.8200 _(0.2572)	0.9067 _(0.1951)	0.9656 _(0.1121)	1.3189	0.0984	0.9256	0.0872	0.1370	0.0872
p = 0.2										
Oracle	—	—	—	—	0.9110	0.1351	0.4973	0.0793	0.2097	0.0793
MLE	—	—	—	—	1.1351	0.1387	0.8888	0.0935	0.1902	0.0935
AdpLASSO	0.6233 _(0.2898)	0.8933 _(0.2459)	0.9483 _(0.1681)	0.9411 _(0.1562)	1.2894	0.1804	0.7608	0.1096	0.2777	0.1096
LASSO	0.6611 _(0.3559)	0.8550 _(0.3030)	0.8750 _(0.2492)	0.7433 _(0.2993)	1.3308	0.2158	0.9016	0.1253	0.2633	0.1253
MCP	0.6400 _(0.2215)	0.9100 _(0.1967)	0.9550 _(0.1433)	0.9589 _(0.1195)	1.0720	0.1371	0.7125	0.0820	0.2793	0.0820
SCAD	0.6478 _(0.2232)	0.9167 _(0.1911)	0.9567 _(0.1409)	0.9567 _(0.1218)	1.0636	0.1386	0.7111	0.0819	0.2787	0.0819
p = 0.3										
Oracle	—	—	—	—	0.8549	0.1253	0.5063	0.0858	0.2964	0.0858
MLE	—	—	—	—	1.0499	0.1552	0.8202	0.0929	0.3081	0.0929
AdpLASSO	0.5878 _(0.2283)	0.8850 _(0.2605)	0.9700 _(0.1189)	0.9322 _(0.1371)	1.1651	0.2074	0.7959	0.1184	0.3692	0.1184
LASSO	0.5822 _(0.3184)	0.8567 _(0.2911)	0.8933 _(0.2246)	0.8422 _(0.2101)	1.2739	0.2632	0.9632	0.1585	0.3676	0.1585
MCP	0.6022 _(0.1645)	0.8900 _(0.2410)	0.9683 _(0.1220)	0.9733 _(0.0906)	0.9990	0.1418	0.6924	0.0799	0.3990	0.0799
SCAD	0.6089 _(0.1624)	0.9017 _(0.2191)	0.9667 _(0.1249)	0.9733 _(0.0946)	0.9956	0.1412	0.6828	0.0823	0.3946	0.0823
p = 0.4										
Oracle	—	—	—	—	0.9225	0.1649	0.5525	0.0980	0.3561	0.0980
MLE	—	—	—	—	1.1974	0.2172	0.9813	0.1080	0.3574	0.1080
AdpLASSO	0.4978 _(0.2993)	0.8600 _(0.2720)	0.9233 _(0.1805)	0.9189 _(0.1581)	1.4347	0.3391	0.9183	0.1653	0.4358	0.1653
LASSO	0.4611 _(0.3782)	0.8283 _(0.3082)	0.8800 _(0.2326)	0.8278 _(0.2305)	1.6320	0.4374	1.1354	0.2329	0.4210	0.2329
MCP	0.5667 _(0.1937)	0.8767 _(0.2380)	0.9250 _(0.1788)	0.9467 _(0.1283)	1.1739	0.2178	0.7907	0.0997	0.4791	0.0997
SCAD	0.5678 _(0.1972)	0.8833 _(0.2307)	0.9233 _(0.1805)	0.9467 _(0.1283)	1.1718	0.2147	0.7915	0.0983	0.4769	0.0983
p = 0.5										
Oracle	—	—	—	—	0.9803	0.1717	0.5442	0.1024	0.4244	0.1024
MLE	—	—	—	—	1.3809	0.2466	1.2807	0.1315	0.4193	0.1315
AdpLASSO	0.4678 _(0.2659)	0.6867 _(0.3660)	0.8450 _(0.2316)	0.9033 _(0.1722)	1.6332	0.3705	1.2317	0.1791	0.5217	0.1791
LASSO	0.4300 _(0.3083)	0.6383 _(0.4007)	0.8383 _(0.2548)	0.8744 _(0.2062)	1.9366	0.4615	1.5444	0.2358	0.5642	0.2358
MCP	0.5144 _(0.2080)	0.7417 _(0.3433)	0.8400 _(0.2336)	0.8956 _(0.1596)	1.4062	0.2975	1.0925	0.1239	0.5012	0.1239
SCAD	0.5156 _(0.2061)	0.7417 _(0.3383)	0.8383 _(0.2343)	0.8900 _(0.1728)	1.4023	0.2993	1.1004	0.1246	0.5083	0.1246

Web Table 6: Results of Example 2: Average (SD) sensitivity, specificity, and estimation errors:
($C = 2, d = 5, n = 400$).

Criteria	Sensitivity		Specificity		$L_2(\hat{\beta}_j)$		$L_2(\hat{\sigma}_j^2)$		$L_2(\hat{\pi}_j)$	
	1	2	1	2	1	2	1	2	1	2
p = 0.05										
Oracle	—	—	—	—	0.9785	0.0572	0.5666	0.0657	0.0609	0.0657
MLE	—	—	—	—	1.3969	0.0650	1.1504	0.0748	0.0758	0.0748
AdpLASSO	0.1967 _(0.2458)	0.7067 _(0.3527)	0.8867 _(0.2213)	0.9789 _(0.1153)	1.6640	0.1458	1.0619	0.1949	0.1355	0.1949
LASSO	0.2322 _(0.3069)	0.6600 _(0.3740)	0.8700 _(0.2582)	0.8533 _(0.2360)	1.6322	0.1719	1.1591	0.2366	0.1429	0.2366
MCP	0.6044 _(0.2404)	0.8433 _(0.2497)	0.8783 _(0.2149)	0.9544 _(0.1240)	1.3575	0.0695	1.0012	0.0881	0.0883	0.0881
SCAD	0.6244 _(0.2339)	0.8583 _(0.2365)	0.8767 _(0.2159)	0.9433 _(0.1394)	1.3490	0.0681	1.0176	0.0861	0.0875	0.0861
p = 0.1										
Oracle	—	—	—	—	0.7820	0.0814	0.4509	0.0739	0.1127	0.0739
MLE	—	—	—	—	1.2497	0.0990	0.9775	0.0916	0.1111	0.0916
AdpLASSO	0.5433 _(0.2936)	0.8217 _(0.2725)	0.9033 _(0.2101)	0.9367 _(0.1638)	1.3122	0.1192	0.8449	0.1247	0.1805	0.1247
LASSO	0.5411 _(0.3670)	0.7817 _(0.3320)	0.8933 _(0.2283)	0.8200 _(0.2442)	1.3590	0.1473	0.9677	0.1489	0.1835	0.1489
MCP	0.6489 _(0.2447)	0.8583 _(0.2330)	0.9083 _(0.1938)	0.9400 _(0.1394)	1.1660	0.0959	0.8370	0.0806	0.1529	0.0806
SCAD	0.6656 _(0.2384)	0.8667 _(0.2252)	0.9100 _(0.1924)	0.9278 _(0.1529)	1.1679	0.0956	0.8470	0.0806	0.1506	0.0806
p = 0.2										
Oracle	—	—	—	—	0.5718	0.0882	0.3571	0.0822	0.1707	0.0822
MLE	—	—	—	—	0.9225	0.1090	0.6945	0.0853	0.1805	0.0853
AdpLASSO	0.6789 _(0.2703)	0.9233 _(0.1939)	0.9650 _(0.1278)	0.9644 _(0.1197)	1.0680	0.1257	0.6026	0.1044	0.2490	0.1044
LASSO	0.4867 _(0.4440)	0.7850 _(0.2856)	0.9367 _(0.1810)	0.9000 _(0.2103)	1.3737	0.2918	0.9929	0.2213	0.2328	0.2213
MCP	0.6644 _(0.1442)	0.9500 _(0.1503)	0.9633 _(0.1306)	0.9733 _(0.0946)	0.9195	0.0933	0.5101	0.0733	0.2413	0.0733
SCAD	0.6667 _(0.1363)	0.9517 _(0.1480)	0.9633 _(0.1306)	0.9733 _(0.0946)	0.9138	0.0927	0.5069	0.0731	0.2403	0.0731
p = 0.3										
Oracle	—	—	—	—	0.5387	0.0911	0.3518	0.0846	0.2383	0.0846
MLE	—	—	—	—	0.7801	0.1099	0.5970	0.0855	0.2577	0.0855
AdpLASSO	0.8533 _(0.2513)	0.9833 _(0.1069)	0.9917 _(0.0641)	0.9922 _(0.0635)	0.8353	0.1166	0.4036	0.0924	0.3013	0.0924
LASSO	0.7033 _(0.4308)	0.9683 _(0.1411)	0.9750 _(0.1092)	0.9711 _(0.1184)	1.1928	0.2378	0.7511	0.1774	0.3229	0.1774
MCP	0.7322 _(0.1583)	0.9900 _(0.0701)	0.9900 _(0.0701)	0.9900 _(0.0631)	0.8046	0.0941	0.3946	0.0786	0.3138	0.0786
SCAD	0.7344 _(0.1573)	0.9900 _(0.0701)	0.9883 _(0.0756)	0.9900 _(0.0631)	0.8013	0.0949	0.3934	0.0787	0.3132	0.0787
p = 0.4										
Oracle	—	—	—	—	0.5475	0.0976	0.3481	0.0840	0.2983	0.0840
MLE	—	—	—	—	1.0006	0.1478	0.8966	0.0999	0.2903	0.0999
AdpLASSO	0.7189 _(0.3605)	0.9233 _(0.2143)	0.9433 _(0.1588)	0.9478 _(0.1654)	1.1094	0.2103	0.7081	0.1482	0.3535	0.1482
LASSO	0.6500 _(0.4008)	0.9167 _(0.2482)	0.9033 _(0.1978)	0.8700 _(0.2209)	1.4012	0.3134	1.0900	0.2231	0.3691	0.2231
MCP	0.7533 _(0.2262)	0.9433 _(0.1588)	0.9367 _(0.1666)	0.9478 _(0.1513)	0.9008	0.1353	0.5396	0.0890	0.3476	0.0890
SCAD	0.7489 _(0.2229)	0.9417 _(0.1608)	0.9367 _(0.1666)	0.9478 _(0.1513)	0.9080	0.1382	0.5417	0.0889	0.3508	0.0889
p = 0.5										
Oracle	—	—	—	—	0.7021	0.1657	0.4427	0.1008	0.3494	0.1008
MLE	—	—	—	—	1.0445	0.1685	0.9230	0.1127	0.3283	0.1127
AdpLASSO	0.5489 _(0.4164)	0.9017 _(0.2794)	0.9433 _(0.1588)	0.9756 _(0.1063)	1.3950	0.3236	0.9470	0.2220	0.4022	0.2220
LASSO	0.5311 _(0.4216)	0.8850 _(0.3156)	0.9100 _(0.1967)	0.8833 _(0.1968)	1.6108	0.3979	1.2915	0.2812	0.4112	0.2812
MCP	0.7178 _(0.1994)	0.9417 _(0.1608)	0.9433 _(0.1588)	0.9322 _(0.1620)	0.9891	0.1540	0.5544	0.0961	0.4057	0.0961
SCAD	0.7078 _(0.1951)	0.9400 _(0.1678)	0.9433 _(0.1588)	0.9433 _(0.1497)	0.9994	0.1551	0.5451	0.0947	0.4111	0.0947

Web Table 7: Average (SD) sensitivity, specificity, and estimation errors over 100 replicates.

$C = 2$		Criteria Component	Sensitivity		Specificity		$L_2(\hat{\beta}_j)$		$L_2(\hat{\sigma}_j^2)$		$L_2(\hat{\pi}_j)$	
d	n		1	2	1	2	1	2	1	2	1	2
Correct model: t FM-VCR												
10	200	Oracle	—	—	—	—	.303	.214	.171	.249	.098	.098
		MLE	—	—	—	—	.474	.420	.333	.326	.153	.153
		AdpLASSO	.482(.235)	.655(.388)	.976(.059)	.968(.062)	.402	.394	.234	.295	.265	.265
		LASSO	.500(.265)	.638(.399)	.955(.072)	.939(.075)	.407	.412	.273	.319	.328	.328
		MCP	.443(.199)	.700(.325)	.993(.038)	.985(.047)	.405	.350	.185	.241	.174	.174
		SCAD	.446(.196)	.702(.320)	.992(.040)	.985(.048)	.404	.346	.186	.235	.167	.167
	400	Oracle	—	—	—	—	.230	.153	.181	.320	.092	.092
		MLE	—	—	—	—	.372	.329	.225	.274	.132	.132
		AdpLASSO	.587(.190)	.940(.195)	.993(.031)	.987(.038)	.327	.228	.180	.190	.172	.172
		LASSO	.556(.250)	.857(.319)	.968(.061)	.953(.063)	.371	.286	.233	.212	.270	.270
		MCP	.556(.175)	.945(.167)	.997(.020)	.995(.024)	.326	.221	.160	.192	.154	.154
		SCAD	.571(.167)	.942(.171)	.997(.020)	.995(.025)	.319	.220	.158	.195	.149	.149
Misspecified model: Gaussian FM-VCR												
d	n											
10	200	Oracle	—	—	—	—	1.083	.342	.786	.181	.330	.181
		MLE	—	—	—	—	1.748	.426	1.667	.234	.443	.234
		AdpLASSO	.394(.272)	.538(.414)	.942(.085)	.939(.076)	1.568	.394	1.489	.309	.387	.309
		LASSO	.391(.335)	.515(.476)	.920(.090)	.925(.086)	1.721	.588	1.721	.425	.558	.425
		MCP	.441(.236)	.588(.368)	.943(.100)	.942(.101)	1.508	.291	1.476	.218	.281	.218
		SCAD	.444(.235)	.583(.370)	.937(.110)	.937(.111)	1.509	.280	1.475	.204	.271	.204
	400	Oracle	—	—	—	—	.809	.282	.553	.111	.199	.111
		MLE	—	—	—	—	1.450	.329	1.210	.209	.325	.209
		AdpLASSO	.544(.247)	.827(.342)	.967(.073)	.970(.061)	1.216	.275	.986	.226	.272	.226
		LASSO	.453(.328)	.650(.467)	.933(.077)	.932(.075)	1.525	.470	1.360	.373	.469	.373
		MCP	.567(.176)	.892(.250)	.957(.113)	.970(.070)	1.249	.201	.897	.177	.219	.177
		SCAD	.573(.173)	.892(.247)	.953(.127)	.970(.076)	1.234	.202	.895	.174	.212	.174

Web Table 8: Average (SD) sensitivity, specificity, and estimation errors for models with orders $C = 1, 2, 3, 4$.

$C = 1$		Criteria Component	Sensitivity		Specificity		$L_2(\hat{\beta}_j)$		$L_2(\hat{\sigma}_j^2)$		$L_2(\hat{\pi}_j)$		
d	n		1	2	1	2	1	2	1	2	1	2	
10	200	Oracle	—	—	—	—	.357	—	.634	—	—	—	
		MLLE	—	—	—	—	.396	—	.616	—	—	—	
		AdpLASSO	.772(.073)	—	.831(.167)	—	.383	—	.672	—	—	—	
		LASSO	.802(.101)	—	.708(.221)	—	.390	—	.673	—	—	—	
		MCP	.771(.069)	—	.859(.151)	—	.375	—	.667	—	—	—	
		SCAD	.767(.064)	—	.874(.137)	—	.371	—	.668	—	—	—	
	400	Oracle	—	—	—	—	.339	—	.558	—	—	—	
		MLLE	—	—	—	—	.365	—	.554	—	—	—	
		AdpLASSO	.984(.102)	—	.712(.243)	—	.365	—	.587	—	—	—	
		LASSO	.972(.079)	—	.738(.255)	—	.375	—	.593	—	—	—	
		MCP	.943(.105)	—	.746(.235)	—	.358	—	.583	—	—	—	
		SCAD	.953(.098)	—	.636(.287)	—	.354	—	.579	—	—	—	
$C = 2$		the true model											
10	200	Oracle	—	—	—	—	.698	.162	.397	.089	.314	.089	
		MLLE	—	—	—	—	1.23	.287	1.02	.138	.335	.138	
		AdpLASSO	.589(.239)	.893(.259)	.970(.067)	.966(.061)	1.09	.288	.848	.161	.501	.161	
		LASSO	.470(.321)	.775(.401)	.928(.100)	.917(.087)	1.32	.383	1.15	.294	.505	.294	
		MCP	.581(.194)	.895(.245)	.984(.045)	.981(.048)	1.04	.238	.763	.106	.503	.106	
		SCAD	.594(.188)	.892(.243)	.983(.049)	.980(.050)	1.03	.237	.761	.107	.501	.107	
	400	Oracle	—	—	—	—	.439	.089	.227	.082	.273	.082	
		MLLE	—	—	—	—	.878	.149	.607	.087	.337	.087	
		AdpLASSO	.690(.127)	.993(.057)	.997(.023)	.995(.029)	.808	.165	.390	.069	.411	.069	
		LASSO	.766(.230)	.970(.166)	.970(.066)	.962(.064)	.862	.232	.518	.108	.425	.108	
		MCP	.692(.104)	.992(.064)	.997(.020)	.998(.014)	.789	.145	.374	.06	.412	.062	
		SCAD	.688(.098)	.992(.064)	.996(.023)	.998(.016)	.792	.146	.377	.062	.412	.062	
	$C = 3$												
	10	200	Oracle	—	—	—	—	.201	.111	.314	.125	.083	.083
			MLLE	—	—	—	—	.371	.297	.173	.115	.216	.092
			AdpLASSO	.654(.140)	.962(.139)	.987(.044)	.992(.033)	.295	.199	.473	.182	.121	.096
			LASSO	.641(.167)	.938(.209)	.967(.065)	.963(.066)	.302	.226	.510	.239	.121	.096
			MCP	.612(.180)	.958(.138)	.997(.025)	.999(.010)	.319	.192	.407	.145	.157	.085
SCAD			.623(.177)	.947(.160)	1.00(.000)	.999(.010)	.310	.192	.412	.153	.150	.084	
400		Oracle	—	—	—	—	.125	.067	.260	.081	.081	.082	
		MLLE	—	—	—	—	.319	.176	.206	.096	.247	.072	
		AdpLASSO	.722(.124)	1.00(.000)	.999(.014)	1.00(.007)	.260	.136	.401	.169	.118	.062	
		LASSO	.717(.119)	.997(.058)	.995(.027)	.994(.027)	.259	.141	.446	.213	.097	.068	
		MCP	.706(.117)	1.00(.000)	1.00(.000)	1.00(.000)	.261	.138	.380	.146	.127	.059	
		SCAD	.708(.116)	1.00(.000)	1.00(.000)	1.00(.000)	.258	.136	.374	.145	.135	.059	
$C = 4$													
10	200	Oracle	—	—	—	—	.199	.106	.313	.113	.083	.083	
		MLLE	—	—	—	—	.389	.322	.119	.098	.294	.165	
		AdpLASSO	.616(.148)	.963(.149)	.994(.028)	.993(.03)	.310	.204	.434	.191	.161	.102	
		LASSO	.642(.140)	.922(.234)	.980(.051)	.977(.05)	.299	.225	.529	.234	.156	.135	
		MCP	.536(.188)	.957(.141)	1.00(.008)	.999(.01)	.372	.205	.342	.133	.237	.090	
		SCAD	.561(.180)	.948(.152)	.999(.012)	.999(.01)	.349	.203	.356	.123	.219	.092	
	400	Oracle	—	—	—	—	.125	.067	.260	.081	.081	.082	
		MLLE	—	—	—	—	.338	.204	.149	.076	.306	.132	
		AdpLASSO	.701(.105)	1.00(.000)	.999(.014)	.999(.010)	.266	.147	.418	.175	.120	.068	
		LASSO	.702(.103)	.998(.029)	.994(.028)	.994(.026)	.265	.150	.452	.226	.104	.076	
		MCP	.693(.116)	1.00(.000)	1.00(.000)	1.00(.000)	.273	.153	.369	.120	.160	.071	
		SCAD	.696(.118)	1.00(.000)	1.00(.000)	1.00(.000)	.268	.154	.356	.115	.176	.075	

Web Table 9: Result of the FMR and FM-VCR models; average (SD) Sensitivity, Specificity, and estimation errors.

$C = 2$		Criteria Component	Sensitivity		Specificity		$L_2(\hat{\beta}_j)$		$L_2(\hat{\sigma}_j)$		$L_2(\hat{\pi}_j)$	
d	n		1	2	1	2	1	2	1	2	1	2
True Gaussian FMR model												
10	200	Oracle	—	—	—	—	.050	.055	.063	.074	.147	.147
		MLLE	—	—	—	—	.115	.140	.138	.158	.232	.232
		AdpLASSO	.529(.364)	.485(.389)	.848(.348)	.837(.347)	.125	.133	.138	.165	.352	.352
		LASSO	.732(.361)	.568(.436)	.843(.279)	.815(.280)	.101	.126	.112	.185	.321	.321
		MCP	.766(.277)	.702(.314)	.864(.249)	.806(.244)	.102	.125	.115	.153	.267	.267
		SCAD	.762(.283)	.725(.287)	.839(.259)	.775(.262)	.106	.125	.120	.145	.259	.259
400	400	Oracle	—	—	—	—	.027	.031	.031	.039	.078	.078
		MLLE	—	—	—	—	.075	.107	.081	.115	.185	.185
		AdpLASSO	.708(.312)	.592(.412)	.878(.318)	.870(.320)	.099	.127	.088	.148	.306	.306
		LASSO	.941(.185)	.790(.374)	.881(.191)	.850(.194)	.060	.087	.053	.118	.209	.209
		MCP	.943(.152)	.888(.242)	.944(.130)	.885(.189)	.050	.075	.062	.094	.169	.169
		SCAD	.951(.144)	.897(.230)	.922(.163)	.865(.222)	.050	.076	.064	.097	.166	.166
Misspecified Gaussian FM-VCR model												
10	200	Oracle	—	—	—	—	.050	.055	.062	.072	.139	.139
		MLLE	—	—	—	—	.099	.132	.060	.098	.133	.133
		AdpLASSO	.562(.323)	.563(.356)	.990(.042)	.970(.060)	.100	.100	.102	.112	.255	.255
		LASSO	.643(.323)	.600(.359)	.991(.034)	.965(.065)	.093	.096	.086	.120	.243	.243
		MCP	.588(.320)	.545(.310)	.991(.040)	.976(.054)	.097	.105	.077	.092	.215	.215
		SCAD	.583(.311)	.548(.304)	.991(.040)	.975(.058)	.097	.100	.066	.080	.202	.202
400	400	Oracle	—	—	—	—	.028	.031	.030	.040	.074	.074
		MLLE	—	—	—	—	.073	.096	.047	.070	.125	.125
		AdpLASSO	.841(.252)	.800(.309)	.993(.032)	.972(.058)	.066	.073	.046	.064	.172	.172
		LASSO	.899(.211)	.810(.325)	.987(.043)	.969(.061)	.059	.072	.048	.073	.177	.177
		MCP	.876(.227)	.810(.299)	.993(.032)	.979(.052)	.054	.062	.039	.055	.144	.144
		SCAD	.872(.234)	.812(.298)	.994(.029)	.981(.051)	.051	.059	.037	.050	.126	.126

Web Table 10: Results of simultaneous selection of (h, λ) with $C = 2$: average (SD) sensitivity, specificity, and estimation errors.

$C = 2$		Criteria Component	Sensitivity		Specificity		$L_2(\hat{\beta}_j)$		$L_2(\hat{\sigma}_j)$		$L_2(\hat{\pi}_j)$	
d	n		1	2	1	2	1	2	1	2	1	2
10	200	Oracle	—	—	—	—	.699	.162	.400	.089	.314	.088
		MLLE	—	—	—	—	1.229	.287	1.015	.136	.335	.136
		AdpLASSO	.604(.234)	.905(.249)	.965(.071)	.962(.066)	1.080	.284	.834	.154	.495	.154
		LASSO	.462(.319)	.783(.396)	.929(.097)	.920(.081)	1.332	.392	1.139	.297	.504	.297
		MCP	.593(.202)	.912(.219)	.981(.052)	.979(.053)	1.035	.236	.743	.106	.497	.106
	SCAD	.608(.194)	.902(.234)	.980(.062)	.975(.054)	1.026	.237	.747	.106	.494	.106	
	400	Oracle	—	—	—	—	.439	.089	.227	.081	.273	.081
		MLLE	—	—	—	—	.878	.149	.606	.086	.337	.086
		AdpLASSO	.768(.167)	.995(.050)	.985(.050)	.993(.037)	.779	.156	.378	.065	.406	.065
		LASSO	.796(.252)	.967(.175)	.946(.088)	.952(.071)	.865	.230	.518	.112	.419	.112
MCP		.766(.162)	.992(.064)	.973(.073)	.997(.019)	.764	.141	.372	.060	.406	.060	
SCAD	.753(.156)	.992(.064)	.963(.089)	.995(.025)	.764	.141	.373	.061	.408	.061		

Web Table 11: Comparison of the methods in terms of tuning parameter and band-width selection time for the adaptive Lasso (AdpLASSO) penalty; $C = 2$;

d	n	Elapsed time	$Pl(h)$ then $BIC(\lambda h)$	Simultaneous $BIC_1(h, \lambda)$
10	200	Time in seconds	3.64	65.45
		Ratio	1	18
	400		9.57	171.17
			1	18

Web Table 12: Average (SD) sensitivity, specificity, and estimation errors for a high-dimensional setting.

$C = 2$	d	n	Criteria Component	Sensitivity		Specificity		$L_2(\hat{\beta}_j)$		$L_2(\hat{\sigma}_j)$		$L_2(\hat{\pi}_j)$	
				1	2	1	2	1	2	1	2	1	2
	500	200	Oracle	—	—	—	—	0.62	0.09	0.30	0.08	0.08	0.08
			MLLE	—	—	—	—	1.71	1.69	0.29	0.35	0.25	0.26
			LASSO	0.50 _(0.31)	0.67 _(0.439)	0.99 _(0.00)	0.99 _(0.00)	1.69	1.59	0.58	0.44	0.54	0.41
			MCP	0.58 _(0.23)	0.80 _(0.28)	0.99 _(0.02)	0.99 _(0.02)	1.55	1.44	0.32	0.31	0.14	0.15
			SCAD	0.55 _(0.23)	0.76 _(0.31)	0.99 _(0.02)	0.99 _(0.02)	1.15	1.41	0.33	0.32	0.16	0.16
		400	Oracle	—	—	—	—	0.35	0.07	0.21	0.07	0.25	0.25
			MLLE	—	—	—	—	2.31	1.94	0.36	0.42	0.36	0.35
			LASSO	0.54 _(0.26)	0.72 _(0.35)	0.99 _(0.00)	0.99 _(0.00)	2.87	2.10	0.34	0.34	0.41	0.41
			MCP	0.61 _(0.25)	0.81 _(0.25)	0.99 _(0.02)	0.99 _(0.02)	2.99	2.20	0.29	0.26	0.40	0.40
			SCAD	0.59 _(0.25)	0.80 _(0.25)	0.99 _(0.01)	0.99 _(0.01)	3.09	2.18	0.29	0.26	0.40	0.39

The runs using adpLASSO did not converge; thus, no results are provided.

Web Table 13: The top 20 covariates with the highest frequencies of being selected based on Pearson’s correlation out of the 100 replications.

	True Covariates				Other Covariates															
x_j	1	2	3	4	74	153	155	5	269	218	268	66	79	329	330	154	118	156	193	103
%	100	100	100	100	99	99	98	96	94	92	92	86	86	80	78	73	70	70	70	62

Web Table 14: The top 20 covariates with the highest frequencies of being selected based on log-likelihood out of the 100 replications.

	True Covariates				Other Covariates															
x_j	1	2	3	4	101	469	352	68	415	118	67	377	74	270	406	414	247	429	5	169
%	100	100	100	100	79	79	78	70	70	65	62	60	56	54	52	52	49	49	48	48