

A New Regression Model: Modal Linear Regression

Weixin Yao and Longhai Li

August 26, 2013

Abstract

The mode of a distribution provides an important summary of data and is often estimated based on some non-parametric kernel density estimator. This article develops a new data analysis tool called modal linear regression in order to explore high-dimensional data. Modal linear regression models the conditional mode of a response Y given a set of predictors \boldsymbol{x} as a linear function of \boldsymbol{x} . Modal linear regression differs from standard linear regression in that standard linear regression models the conditional mean (as opposed to mode) of Y as a linear function of \boldsymbol{x} . We propose an Expectation-Maximization algorithm in order to estimate the regression coefficients of modal linear regression. We also provide asymptotic properties for the proposed estimator using the algorithm under the assumption of a skewed error density. Our empirical studies with simulated data and real data demonstrate that the proposed modal regression gives shorter predictive intervals than mean linear regression, median linear regression, and MM-estimators.

Key Words: Forest fire data; Linear regression; Modal regression; Mode.

¹Weixin Yao is Assistant Professor, Department of Statistics, Kansas State University, Manhattan, Kansas 66506, U.S.A. Email: wxyao@ksu.edu. Longhai Li is Associate Professor, Department of Mathematics and Statistics, University of Saskatchewan, Saskatoon, Saskatchewan, S7N5E6, Canada. Email: longhai@math.usask.ca. The research of Longhai Li is supported by fundings from Natural Sciences and Engineering Research Council of Canada, and Canadian Foundation of Innovations.

1. INTRODUCTION

Mode provides an important summary of data. Many authors have made efforts to identify modes of population distributions for low-dimensional data (see for example Muller and Sawitzki (1991); Scott (1992); Friedman and Fisher (1999); Chaudhuri and Marron (1999); Fisher and Marron (2001); Davies and Kovac (2004); Hall, Minnotte, and Zhang (2004); Ray and Lindsay (2005); Yao and Lindsay (2009), as well as documentations of the R package “np” for non-parametric mode estimation). In high-dimensional data, it is often of interest to estimate conditional distributions in order to identify associations between a response and a set of predictors. To the best of our knowledge, little research has been done to hunt conditional modes in regression problems.

Suppose we have collected a random sample $\{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$, where \mathbf{x}_i is a p -dimensional column vector and y_i is observation of a continuous response variable Y . We are interested in estimating the conditional density function of Y given \mathbf{x} , denoted by $f(y | \mathbf{x})$. Conventional regression methods usually model the mean or median of $f(y | \mathbf{x})$ as a linear function of \mathbf{x} . Another regression approach is to model the *mode* of $f(y | \mathbf{x})$ as a linear function of \mathbf{x} (Lee, 1989). We will refer to this approach as *modal linear regression* (or MODLR for short) throughout this paper. Compared to other regression approaches, modal linear regression has the following features:

1. MODLR attempts to capture the “most probable” value — the mode (instead of the mean, median, or quantile) of the conditional distribution of Y given \mathbf{x} . The conditional mode may be a more useful summary than the conditional mean, median, or quantile when the conditional distribution of Y given \mathbf{x} is asymmetric.
2. MODLR may provide shorter prediction intervals than other linear regression approaches for a nominal confidence level, since an interval around a conditional mode can cover more samples than an interval of the same length around a conditional mean.

3. MODLR is robust to outliers that don't follow the same relationship exhibited by the majority of a sample and is particularly robust to heavy-tailed conditional error distributions. This is because modal regression focuses on modelling the majority of the distribution of Y given \mathbf{x} .
4. MODLR is well justified in situations where conditional distributions are highly skewed. Many robust regression methods, such as median regression and MM-estimators, require symmetries in conditional distributions in order to achieve good performance. Quantile regression methods allow for skewed conditional distributions, but MODLR gives estimates of conditional modes which may be of more interest than a quantile in some application contexts.

Modal linear regression is potentially a very useful addition to current data analysis tools. However, estimation of modal regression coefficients is not trivial. In this article, we propose an EM algorithm that minimizes a kernel-based objective function for estimating modal regression coefficients. We have studied asymptotic and other theoretical properties of the proposed estimation procedure. We also propose a method for constructing asymmetric prediction intervals that can have better coverage than symmetric prediction intervals when conditional distributions are highly skewed.

The rest of this article is organized as follows. In Section 2, we introduce the kernel-based objective function and the EM algorithm for maximizing it; we also provide the theoretical properties of the estimating procedure. In Section 3, we use simulated datasets to compare the proposed MODLR with least square regression, median regression, and MM-estimators. We also compare these regression methods using forest fire data. Our empirical results show that MODLR provides significantly shorter prediction intervals than other regression methods. The article is concluded in Section 4 with discussions of possible future work. Proofs of the consistency of our estimators and necessary technical conditions are given in

the Appendix.

2. MODAL LINEAR REGRESSION

2.1. Introduction to modal linear regression

Suppose that a response variable Y given a set of predictor \mathbf{x} is distributed with a probability density function $f(y | \mathbf{x})$. Assume that the mode of $f(y | \mathbf{x})$, denoted by $\text{Mode}(Y | \mathbf{x}) = \arg \max_y (f(y | \mathbf{x}))$, is unique. Furthermore, assume that $\text{Mode}(Y | \mathbf{x})$ is a linear function of \mathbf{x} , *i.e.*

$$\text{Mode}(Y | \mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}. \quad (2.1)$$

In (2.1) we assume that the first element of \mathbf{x} is 1; this represents the intercept term. Let $\epsilon = y - \mathbf{x}^T \boldsymbol{\beta}$; we denote the conditional density of ϵ given \mathbf{x} by $g(\epsilon | \mathbf{x})$ and refer to it as the *error distribution*. Note that the estimation method (and its asymptotic justification) that we will propose next allows for the error distribution to depend on \mathbf{x} .

If $g(\epsilon | \mathbf{x})$ is symmetric about 0, the $\boldsymbol{\beta}$ in (2.1) will be the same as the coefficients obtained by conventional linear regression; however, if $g(\epsilon | \mathbf{x})$ is skewed, modal regression coefficients and conventional linear regression coefficients will be different. It is possible that the mode of Y given \mathbf{x} is a linear function of \mathbf{x} but the conventional mean is nonlinear. The following example illustrates the difference between modal regression function and conventional mean regression function when the error distribution is skewed.

Example 1: Let (\mathbf{x}, Y) satisfy the following model assumption

$$Y = m(\mathbf{x}) + \sigma(\mathbf{x})\epsilon, \quad (2.2)$$

where ϵ has density $h(\cdot)$. Suppose $h(\cdot)$ is a skewed density with mean 0 and mode 1.

a) If $m(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}$ and $\sigma(x) = \mathbf{x}^T \boldsymbol{\alpha}$, then

$$E(Y | \mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta} \quad \text{and} \quad \text{Mode}(Y | \mathbf{x}) = \mathbf{x}^T (\boldsymbol{\beta} + \boldsymbol{\alpha}).$$

Thus, Y depends on \mathbf{x} linearly from the point of view of both mean regression and modal regression even though their regression parameters are different.

b) If $m(\mathbf{x}) = 0$ and $\sigma(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\alpha}$, then

$$E(Y | \mathbf{x}) = 0 \quad \text{and} \quad \text{Mode}(Y | \mathbf{x}) = \mathbf{x}^T \boldsymbol{\alpha}.$$

Therefore, in terms of conditional mean, Y does not depend on \mathbf{x} ; however, in terms of conditional mode, Y does depend linearly on \mathbf{x} . From this example we see that variable selection techniques based on modal regression might reveal some useful predictors when mean regression cannot.

To estimate the modal regression parameter $\boldsymbol{\beta}$ in (2.1), we propose **maximizing** the kernel-based objective function

$$Q_h(\boldsymbol{\beta}) \equiv \frac{1}{n} \sum_{i=1}^n \phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}), \quad (2.3)$$

where $\phi_h(t) = h^{-1} \phi(t/h)$ and $\phi(t)$ is a kernel density function symmetric about 0. For the remainder of the paper we will assume that ϕ is the standard normal density (for simplicity of computation). Based on this choice of kernel, the M-step of the MEM algorithm presented next has the closed-form solution shown in Equation (2.6). It should be noted that all the asymptotic results presented in this article still hold if other kernels are used. We will denote the maximizer of (2.3) by $\hat{\boldsymbol{\beta}}$ and call it the modal linear regression estimator, shortened by MODLRE.

We now explain why (2.3) can be used to estimate the modal regression coefficients. We first look at the simplest case in which there is no predictor, *i.e.* $\beta = \beta_0$. For such cases, (2.3) is simplified to become

$$Q_h(\beta_0) \equiv \frac{1}{n} \sum_{i=1}^n \phi_h(y_i - \beta_0), \quad (2.4)$$

Where $Q_h(\cdot)$ is the kernel estimate of the density function of Y . Therefore, the maximizer of (2.4) is the mode of the kernel density function based on y_1, \dots, y_n . When $n \rightarrow \infty$ and $h \rightarrow 0$, the mode of this kernel density function will converge to the mode of the distribution of Y . Such a modal estimator has been proposed by Parzen (1962). When there are predictors present, for any fixed β , $Q_h(\beta)$ in (2.3) is the value of the kernel density function based on the residuals $\epsilon_i = y_i - x_i\beta$ at $\epsilon = 0$. Maximizing (2.3) with respect to β yields the line $x\hat{\beta}$ such that the kernel density function of residuals ϵ_i has highest value at 0. In the special case that $\phi_h(t) = (2h)^{-1}I(|t| \leq h)$, a uniform kernel, maximizing (2.3) yields the line $x^T\hat{\beta}$ such that the band $x^T\hat{\beta} \pm h$ covers the largest number of response values y_i .

Lee (1989) used a uniform kernel to estimate modal regression coefficients. In his theoretical investigation, h is fixed and does not depend on the sample size n . In order to get consistency results for the estimator, Lee assumed the error distribution to be symmetric. Note that in such cases the modal line is the same as the traditional mean regression line. Thus, Lee's theoretical results didn't justify applications of MODLR for situations with skewed error distributions (where MODLR is more useful than other regression methods). In this article, we prove (see Appendix for details) that if we let $h \rightarrow 0$ when $n \rightarrow \infty$, the $\hat{\beta}$ found by maximizing $Q_h(\beta)$ in (2.3) is a consistent estimate of the modal regression parameter in (2.1) for very general error density functions without symmetry assumptions.

2.2. Modal EM algorithm

There is no closed-form expression of the maximizer of (2.3); therefore, we propose to

extend the modal expectation-maximization (MEM) algorithm (Li, Ray, and Lindsay, 2007; Yao, 2013) in order to maximize (2.3).

Similar to an EM algorithm, the MEM algorithm consists of an E-step and an M-step: Starting with $\boldsymbol{\beta}^{(0)}$, repeat the following two steps until it converges:

E-Step: In this step, we calculate weights $\pi(j | \boldsymbol{\beta}^{(k)})$, $j = 1, \dots, n$ as

$$\pi(j | \boldsymbol{\beta}^{(k)}) = \frac{\phi_h(y_j - \mathbf{x}_j^T \boldsymbol{\beta}^{(k)})}{\sum_{i=1}^n \phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k)})} \propto \phi_h(y_j - \mathbf{x}_j^T \boldsymbol{\beta}^{(k)}). \quad (2.5)$$

M-Step: In this step, we update $\boldsymbol{\beta}^{(k+1)}$

$$\begin{aligned} \boldsymbol{\beta}^{(k+1)} &= \arg \max_{\boldsymbol{\beta}} \sum_{j=1}^n \left\{ \pi(j | \boldsymbol{\beta}^{(k)}) \log \phi_h(y_j - \mathbf{x}_j^T \boldsymbol{\beta}) \right\} \\ &= (\mathbf{X}^T \mathbf{W}_k \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_k \mathbf{y}, \end{aligned} \quad (2.6)$$

where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$, \mathbf{W}_k is an $n \times n$ diagonal matrix with diagonal elements $\pi(j | \boldsymbol{\beta}^{(k)})$ s, and $\mathbf{y} = (y_1, \dots, y_n)^T$.

Some remarks on the proposed MEM algorithm:

1. The major difference between the least squares estimate (LSE) and the modal regression estimate (MODLRE) lies in the E step. For the LSE, each observation has equal weights, whereas for MODLRE the weights depend on how close y_i is to the modal regression line. This weighting scheme allows MODLRE to reduce the effect of observations far away from the modal regression line in order to achieve robustness.
2. When the normal kernel is used for ϕ in (2.3), the function optimized in the M-step is a weighted sum of log likelihoods corresponding to ordinary linear regression. In this case we obtain a closed-form expression for the maximizer in (2.6). If other kernels are used, some optimization algorithms are needed in the M-step.

3. The converged value obtained by the MEM algorithm depends on the starting point chosen, and there is no guarantee that the algorithm will converge to the global optimal solution of (2.3). Therefore, it is prudent to run the algorithm multiple times using several different starting points and choose the best local optima found.

We have proved (see Appendix) the ascending property of the proposed MEM for any choice of kernel for ϕ in (2.3):

Theorem 2.1. *Each iteration of (2.5) and (2.6) will monotonically non-decrease the objective function (2.3), i.e., $Q_h(\boldsymbol{\beta}^{(k+1)}) \geq Q_h(\boldsymbol{\beta}^{(k)})$, for all k .*

The iteratively reweighted least squares (IRWLS) algorithm has been commonly used for general M-estimators. Since the maximizer of (2.3) can be considered as a special case of M-estimators, the IRWLS algorithm can be applied to find $\hat{\boldsymbol{\beta}}$. When the normal kernel $\phi(\cdot)$ is used, the IRWLS algorithm is indeed equivalent to the proposed MEM algorithm, but when other kernels are used, the two algorithms are different. IRWLS has been proven to be ascending (*i.e.* monotonically non-decreases the objective function) if $-\phi(x)/x$ is non-increasing (Huber, 1981). However, when $\phi(x)$ is a normal density function, $-\phi(x)/x$ is not non-increasing. Therefore the existing theories of IRWLS cannot justify Theorem 2.1 if the normal kernel $\phi(\cdot)$ is used. Because the proof of Theorem 2.1 is for any kernel density $\phi(\cdot)$, including the normal kernel, Theorem 2.1 provides an extension to existing IRWLS theories.

2.3. Asymptotic properties of $\hat{\boldsymbol{\beta}}$

The asymptotic properties established for traditional M-estimators are based on assumptions that the error density is symmetric and the objective function is fixed. In addition, the target of traditional M-estimators is conditional mean. For our proposed modal regression, we will allow that the tuning parameter h in the objective function goes to zero and the error density can be skewed. Therefore, the theoretical results on the traditional M-estimators

cannot be directly applied to the proposed modal linear regression estimator. In this section, we will give the results about the consistency of our proposed modal regression estimator $\hat{\boldsymbol{\beta}}$ for model (2.1), its convergence rate, and its asymptotic distribution. Their proofs are given in the Appendix.

Theorem 2.2. *When $h \rightarrow 0$ and $nh^5 \rightarrow \infty$, under the regularity conditions (A1)–(A3) in the Appendix, there exists a consistent maximizer of (2.3) such that*

$$\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\| = O_p\{h^2 + (nh^3)^{-1/2}\},$$

where $\boldsymbol{\beta}_0$ is the true coefficient of the modal regression function defined in (2.1).

Theorem 2.3. *Under the same assumptions as Theorem 2.2, the $\hat{\boldsymbol{\beta}}$ that satisfies the consistency result given in Theorem 2.2, has the following asymptotic normality result*

$$\sqrt{nh^3} \left[\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 - \frac{h^2}{2} J^{-1} K \{1 + o_p(1)\} \right] \xrightarrow{D} N \{0, \nu_2 J^{-1} L J^{-1}\}, \quad (2.7)$$

where $\nu_2 = \int t^2 \phi^2(t) dt$ and

$$J = E\{g''(0 | \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T\}; K = E\{g'''(0 | \mathbf{x}_i) \mathbf{x}_i\}; L = E\{g(0 | \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T\}.$$

Pazen (1962) and Eddy (1980) have proven similar asymptotic results for kernel estimators of the mode of the distribution of Y without conditioning on \mathbf{x} . Therefore, the results of Pazen (1962) and Eddy (1980) can be considered as a special case of Theorem 2.3 when there is no predictor involved, i.e., $\mathbf{x} = 1$.

By Theorem 2.3, the asymptotic bias of $\hat{\boldsymbol{\beta}}$ is $h^2 J^{-1} K / 2$ and the asymptotic variance is $\nu_2 J^{-1} L J^{-1} / (nh^3)$. A theoretic optimal bandwidth h for estimating $\boldsymbol{\beta}$ can be obtained by

minimizing the asymptotic weighted mean squared errors (MSE)

$$\mathbb{E} \left\{ (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^T W (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \right\} \approx K^T J^{-1} W J^{-1} K h^4 / 4 + (nh^3)^{-1} \nu_2 \text{tr} (J^{-1} L J^{-1} W),$$

where $\text{tr}(A)$ is the trace of A and W is a diagonal matrix, whose diagonal elements reflect the importance of the accuracy in estimating different coefficients. Therefore, the asymptotic optimal bandwidth h is

$$\hat{h}_{opt} = \left[\frac{3\nu_2 \text{tr} (J^{-1} L J^{-1} W)}{K^T J^{-1} W J^{-1} K} \right]^{1/7} n^{-1/7}. \quad (2.8)$$

If $W = (J^{-1} L J^{-1})^{-1} = J L^{-1} J$, which is proportional to the inverse of the asymptotic variance of $\hat{\boldsymbol{\beta}}$, then

$$\hat{h}_{opt} = \left[\frac{3\nu_2(p+1)}{K^T L^{-1} K} \right]^{1/7} n^{-1/7}. \quad (2.9)$$

Let $\boldsymbol{\beta} = (\beta_0, \boldsymbol{\beta}_s)^T$, where β_0 is a scalar intercept parameter and $\boldsymbol{\beta}_s$ is the slope parameter.

If ϵ is independent of \boldsymbol{x} , then

$$J^{-1} K = (1, 0, \dots, 0)^T g'''(0) / \{2g''(0)\},$$

and thus the asymptotic bias of the slope parameter $\boldsymbol{\beta}_s$ is 0. Therefore, the optimal bandwidth h for estimating $\boldsymbol{\beta}_s$ should go to infinity, which implies that the resulting estimate $\hat{\boldsymbol{\beta}}_s$ is a least square estimate with root n consistency. This is expected since when ϵ is independent of \boldsymbol{x} , the slope parameter $\boldsymbol{\beta}_s$ of modal regression line is the same as the slope parameter of conventional mean regression line and thus can be estimated at root n convergence rate.

Given the root n consistent estimate $\hat{\boldsymbol{\beta}}_s$ (using LSE, for example), we propose to further estimate β_0 by

$$\hat{\beta}_0 = \arg \max_{\beta_0} \frac{1}{n} \sum_{i=1}^n \phi_h(y_i - x_i^T \hat{\boldsymbol{\beta}}_s - \beta_0). \quad (2.10)$$

The above maximization can be done similarly using the MEM algorithm proposed in Section 2.2. We have the following result for $\hat{\beta}_0$. Its proof is given in the Appendix.

Theorem 2.4. *Under the same assumption as Theorem 2.2, if ϵ is independent of \mathbf{x} and $g''(0) \neq 0$, the $\hat{\beta}_0$ defined in (2.10) has the following asymptotic distribution:*

$$\sqrt{nh^3} \left\{ \hat{\beta}_0 - \beta_0 - \frac{g'''(0)h^2}{2g''(0)} + o_p(h^2) \right\} \xrightarrow{D} N \left\{ 0, \frac{g(0)\nu_2}{[g''(0)]^2} \right\}. \quad (2.11)$$

Note that when ϵ is independent of \mathbf{x} , $J^{-1}LJ^{-1} = g''(0)^{-2}g(0)E(\mathbf{x}\mathbf{x}^T)^{-1}$. Let

$$A = E(\mathbf{x}\mathbf{x}^T) = \begin{pmatrix} 1 & A_{12} \\ A_{12}^T & A_{22} \end{pmatrix}$$

and a^{11} be the (1,1) element of A^{-1} . Noting that $a^{11} = (1 - A_{12}A_{22}^{-1}A_{12}^T)^{-1}$ and A_{22} is positive definite, we have $a^{11} \geq 1$. Therefore, based on Theorem 2.3 and 2.4, we can see that using the root n consistent estimate $\hat{\beta}_s$ as initial, we can get more efficient estimate of the intercept β_0 than the one found by maximizing (2.3) directly. This is reasonable since the estimate $\hat{\beta}_0$ in (2.10) need not account for the uncertainty of $\hat{\beta}_s$ due to its root n consistency and thus $\hat{\beta}_0$ is asymptotically as efficient as if β_s were known.

From Theorem 2.4, we can see that the asymptotic bias of $\hat{\beta}_0$ is $\{2g''(0)\}^{-1}g'''(0)h^2$ and its asymptotic variance is $[\{g''(0)\}^2nh^3]^{-1}g(0)\nu_2$. By minimizing the asymptotic MSE, we can get the asymptotic optimal bandwidth h for estimating β_0 :

$$\hat{h}_{opt} = \left[\frac{3g(0)\nu_2}{\{g'''(0)\}^2} \right]^{1/7} n^{-1/7}. \quad (2.12)$$

2.4. Finite sample breakdown point

To investigate robustness of the MODLRE, we also calculate its finite sample breakdown point. A breakdown point is used to quantify the proportion of bad data in a sample that an

estimator can tolerate before returning arbitrary values. Since usually the breakdown point is most useful in a small sample setup (Donoho, 1982; Donoho and Huber, 1983), we will mainly focus on the finite sample breakdown point. A number of definitions for the finite sample breakdown point have been proposed (see, for example, Hampel, 1971, 1974; Donoho, 1982; Donoho and Huber, 1983). In this paper, we shall work with the finite sample contamination breakdown point. Let $\mathbf{z}_i = (\mathbf{x}_i, y_i)$. Given the sample $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)$, denote $T(\mathbf{Z})$ the MODLRE $\hat{\boldsymbol{\beta}}$, as defined as the maximizer of (2.3). We can corrupt the original sample \mathbf{Z} by adding m arbitrary points $\mathbf{Z}' = (\mathbf{z}_{n+1}, \dots, \mathbf{z}_{n+m})$. The corrupted sample $\mathbf{Z} \cup \mathbf{Z}'$ then has sample size $n + m$ and contains a fraction $\delta = m/(m + n)$ of bad values. The finite sample contamination breakdown point δ^* is defined as

$$\delta^*(\mathbf{Z}, T) = \min_{1 \leq m \leq n} \left\{ \frac{m}{n + m} : \sup_{\mathbf{Z}'} \|T(\mathbf{Z} \cup \mathbf{Z}')\| = \infty \right\}, \quad (2.13)$$

where $\|\cdot\|$ is Euclidean norm.

Theorem 2.5. *Given observations $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)$, suppose $T(\mathbf{Z}) = \hat{\boldsymbol{\beta}}$, the MODLRE estimate defined as the maximizer of (2.3). Let*

$$M = \sqrt{2\pi}h \sum_{i=1}^n \phi_h(y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}). \quad (2.14)$$

Then the finite sample contamination breakdown point of MODLRE is

$$\delta^*(\mathbf{Z}, T) = \frac{m^*}{n + m^*}, \quad (2.15)$$

where m^ is an integer satisfying $\lceil M \rceil \leq m^* \leq \lfloor M \rfloor + 1$, $\lfloor a \rfloor$ is the largest integer not greater than a , and $\lceil a \rceil$ is the smallest integer not less than a .*

The proof of Theorem 2.5 is given in the Appendix. From the above theorem, we can see

that the breakdown point depends not only on $\phi(\cdot)$, and the tuning parameter h , but also on the sample configuration. (However, Huber (1984) pointed out if the scale (contained in the bandwidth h of the MODLRE) is determined from the sample itself, empirically, the breakdown point is quite high.)

3. SIMULATION STUDY AND APPLICATION

In this section we conduct a Monte Carlo simulation study in order to assess the performance of our proposed MODLRE under a finite sample size scenario. We will compare MODLRE with some other regression methods. A real data application is also provided.

3.1. Bandwidth selection

The modal regression estimator requires a selection of the bandwidth. The asymptotically optimal bandwidth formula (2.9) contains the unknown quantities $g^{(v)}(0 | \mathbf{x})$, $v = 0, 2, 3$, *i.e.* the v th derivative of the conditional density of ϵ given \mathbf{x} . Hence, they are not ready to use. A commonly used method is to replace these unknown quantities with estimates. Given the initial residual $\hat{\epsilon}_i = y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}$, where $\hat{\boldsymbol{\beta}}$ is the traditional least squares estimate (or a robust estimate if there are some outliers) of $\boldsymbol{\beta}$, we can estimate their mode, denoted by \hat{m} , by maximizing the kernel density estimator (Pare, 1962). Under the assumption of independence of ϵ and \mathbf{x} , $\hat{\epsilon}_i - \hat{m}$ approximately has density $g(\cdot)$ and thus $g^{(v)}(0 | \mathbf{x})$ can be estimated by (see, for example, Silverman, 1986 and Scott, 1992)

$$\hat{g}^{(v)}(0 | \mathbf{x}) = \frac{1}{nh^{v+1}} \sum_{i=1}^n K^{(v)} \left\{ \frac{\hat{\epsilon}_i - \hat{m}}{h} \right\}, \quad v = 0, 2, 3,$$

where h is chosen using the method reported by Botev et.al. (2010) and $K^{(v)}(\cdot)$ is the v th derivative of kernel density function $K(\cdot)$. Then we can estimate J , K , and L by

$$\hat{J} = n^{-1} \sum_{i=1}^n \hat{g}''(0 | \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T, \hat{K} = n^{-1} \sum_{i=1}^n \hat{g}'''(0 | \mathbf{x}_i) \mathbf{x}_i, \text{ and } \hat{L} = n^{-1} \sum_{i=1}^n \hat{g}(0 | \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T,$$

and apply Equation (2.9) to estimate \hat{h}_{opt} . To refine the bandwidth selection, one might further iteratively update a chosen bandwidth by recalculating the residual $\hat{\epsilon}_i$ given by the modal linear regression estimate.

3.2. A Monte Carlo simulation study

We generated an iid sample $\{(x_i, y_i), i = 1, \dots, n\}$ from the following model

$$Y = 1 + 3X + \sigma(X)\varepsilon,$$

where $X \sim U(0, 1)$, $\varepsilon \sim 0.5N(-1, 2.5^2) + 0.5N(1, 0.5^2)$, and $\sigma(X) = 1 + 2X$. Note that $E(\varepsilon) = 0$, $\text{Mode}(\varepsilon) = 1$, and $\text{Median}(\varepsilon) = 0.67$ (the last two quantities are approximate). For this model, the conditional mean regression function is $E(Y | X) = 1 + 3X$ and the conditional modal regression function is $\text{Mode}(Y | X) = 2 + 5X$. The modal regression residual is $Y - \text{Mode}(Y|X) = (1 + 2X)(\varepsilon - 1)$, whose distribution peaks at 0 but is negatively skewed. The conditional median function is $\text{Median}(Y|X) = 1.67 + 4.34X$. We consider and compare the following four methods: 1) traditional mean regression based on the least squares estimator (LSE); 2) median regression (MEDREG); 3) MM-estimate (M-estimate with a initial robust M-estimate of scale, Yohai, 1987) based on Tukey's biweight ψ -function; 4) the proposed modal linear regression (MODLR).

Figure 1 shows the scatter plot of a typical generated sample with $n = 200$, as well as regression lines corresponding to the 4 regression methods. From the plot we can see that the modal regression line goes through the area containing the most number of points. A small prediction band around this line is expected to contain the most number of future points. In contrast, the mean regression line based on LSE is *skewed* to a flatter line and lies in a much less dense area for capturing the conditional mean. The regression lines based on the median regression and the MM-estimate lie in higher density areas than the regression line based on LSE.

Table 1 reports the average and standard error (Std) of the parameter estimates for each method based on 1,000 replicates. From this table we see that LSE, MEDREG, and MODLRE estimate their target parameters well. However, the MM-estimate does not estimate the conditional mean function well; this is because the assumption of symmetric error density is violated. Surprisingly, the MODLRE has smaller standard error than the other methods in this example (when the error is skewed), especially when $n = 200$ or $n = 400$. Therefore for finite samples the MODLRE not only has a good modal explanation but also might have better estimation accuracy than other methods when the error is skewed.

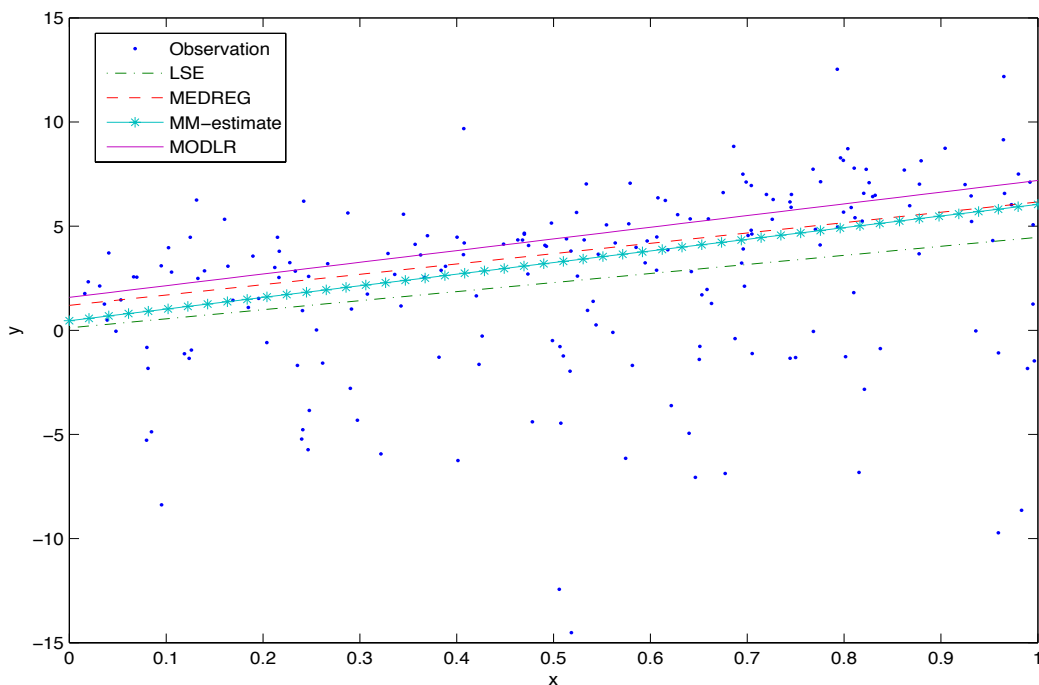


Figure 1: Scatter plot of a typical sample with $n = 200$ for Example 1 with different estimated regression lines: ‘-.’ denotes the mean regression line based on LSE; ‘--’ denotes the median regression line; ‘-*’ denotes the regression line based on MM-estimate; ‘-’ denotes the modal regression line.

Table 2 reports the average (and standard error) of the coverage probabilities of prediction intervals of similar lengths centered around each estimated regression line in 1,000 replicates.

We consider three different lengths of intervals: 0.1σ , 0.2σ , and 0.5σ , where $\sigma = 2$ is the approximate standard error of ε . For each of the 1,000 replications the coverage probability is estimated from 1,000 new cases where the predictor \mathbf{x} is equally spaced between 0.1 to 0.9. From Table 2 we see that MODLRE provides higher coverage probabilities than the other three methods. In addition, MEDREG provides larger coverage probabilities than the MM-estimate and LSE, while the MM-estimate provides larger coverage probabilities than LSE. Note that when the lengths of these intervals are large enough the different methods will provide similar coverage probabilities.

Table 1: Average (Std) of parameter estimates over 1,000 repetitions.

Method	Parameter	n=50	n=100	n=200	n=400
LSE	$\beta_0 = 1$	1.022(0.964)	0.989(0.659)	1.007(0.490)	1.009(0.322)
	$\beta_1 = 3$	2.890(2.260)	3.063(1.500)	2.977(1.160)	2.976(0.733)
MEDREG	$\beta_0 = 1.67$	1.587(0.707)	1.613(0.422)	1.636(0.301)	1.667(0.188)
	$\beta_1 = 4.34$	4.226(1.670)	4.372(0.981)	4.339(0.705)	4.312(0.457)
MM-estimate	$\beta_0 = 1$	1.051(0.782)	1.040(0.530)	1.022(0.376)	1.035(0.265)
	$\beta_1 = 3$	5.123(1.640)	5.234(1.060)	5.271(0.744)	5.271(0.512)
MODLRE	$\beta_0 = 2$	1.789(0.670)	1.841(0.372)	1.875(0.229)	1.912(0.140)
	$\beta_1 = 5$	4.829(1.750)	5.024(0.948)	5.044(0.574)	5.020(0.387)

Table 2: Average (Std) of coverage probabilities over 1,000 repetitions with $\sigma = 2$.

Width	Method	n=50	n=100	n=200	n=400
0.1 σ	LSE	0.034(0.015)	0.032(0.011)	0.030(0.009)	0.029(0.007)
	MEDREG	0.073(0.018)	0.077(0.014)	0.078(0.012)	0.080(0.010)
	MM-estimate	0.065(0.023)	0.067(0.019)	0.066(0.015)	0.067(0.012)
	MODLRE	0.087(0.016)	0.092(0.012)	0.095(0.010)	0.095(0.009)
0.2 σ	LSE	0.069(0.028)	0.065(0.022)	0.061(0.015)	0.059(0.013)
	MEDREG	0.144(0.033)	0.153(0.024)	0.155(0.019)	0.158(0.015)
	MM-estimate	0.129(0.042)	0.133(0.034)	0.132(0.027)	0.134(0.021)
	MODLRE	0.170(0.027)	0.179(0.018)	0.184(0.013)	0.186(0.012)
0.5 σ	LSE	0.186(0.062)	0.181(0.047)	0.174(0.035)	0.171(0.028)
	MEDREG	0.338(0.061)	0.355(0.040)	0.360(0.029)	0.365(0.022)
	MM-estimate	0.313(0.080)	0.322(0.062)	0.325(0.046)	0.330(0.036)
	MODLRE	0.378(0.049)	0.395(0.029)	0.404(0.018)	0.407(0.015)

3.3. Application to forest fire data

Forest fires, also called wildfires, cause great ecological and economical damage. Fast detection of a forest fire is vital for successful fire fighting, but traditional human or automatic surveillance (such as by satellites, infrared or smoke scanners) is expensive. Recently the use of low-cost meteorological data (such as temperature, wind, and precipitation data) to warn the public of a potential wildfire has received a lot of attention. This inexpensive form of information can also be used to get a quick estimate of post-fire damage.

In this section we compare the proposed MODLR and other regression techniques with a forest fire dataset (Cortez and Morais, 2007). The data was downloaded from <http://www.dsi.uminho.pt/~pcortez/forestfires>. This forest fire data contains 517 observations and was collected between January 2000 and December 2003 from the Montesinho natural park of the Trás-os-Montes northeast region of Portugal. On a daily basis, every time a forest fire occurred many features were recorded, such as the time, date, spatial location, and weather conditions. Following Cortez and Morais (2007), we use four meteorological variables: outside temperature (temp), outside relative humidity (RH), outside wind speed

(wind), and outside rain (rain), as predictors for the total burned area (area). We fit the data by LSE, MEDREG, MM-estimate, and MODLRE. One important feature of this dataset is that it contains outliers and a positively-skewed response variable (area); therefore it is expected that the the proposed modal linear regression will compare favorably to other methods.

To compare the four regression methods we look at the widths of each prediction interval (with the same confidence level). For constructing confidence intervals, we assume that the error distribution of ϵ is independent of \mathbf{x} . Suppose we have obtained the parameter estimate $\hat{\boldsymbol{\beta}}$ and the corresponding error (residual) $\hat{\epsilon}_i = y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}$ for $i = 1, \dots, n$; we will use $\hat{\epsilon}_{[i]}$ to denote the i th smallest value of the residuals. The traditional prediction interval with confidence level α for the new predictor \mathbf{x}_{new} is symmetric about the point prediction of y_{new} : $(\mathbf{x}_{new}^T \hat{\boldsymbol{\beta}} - \hat{\epsilon}_{[n_1]}, \mathbf{x}_{new}^T \hat{\boldsymbol{\beta}} + \hat{\epsilon}_{[n_2]})$, where $n_1 = \lfloor n\alpha/2 \rfloor$ and $n_2 = n - n_1$. This symmetric method will be ideal if the regression error distribution is symmetric. To consider and make use of the skewness of the error distribution, we propose to construct **asymmetric prediction intervals** as follows. Suppose $\hat{g}(\cdot)$ is the kernel density estimate of ϵ based on the residuals $\hat{\epsilon}_1, \dots, \hat{\epsilon}_n$ that are estimated by MODLRE. We propose to find the indexes $k_1 < k_2$ such that $k_2 - k_1 = n_2 - n_1 = \lceil n(1 - \alpha) \rceil$ and $\hat{g}(\hat{\epsilon}_{[k_1]}) \approx \hat{g}(\hat{\epsilon}_{[k_2]})$. The proposed prediction interval for the new predictor \mathbf{x}_{new} is $(\mathbf{x}_{new}^T \hat{\boldsymbol{\beta}} - \hat{\epsilon}_{[k_1]}, \mathbf{x}_{new}^T \hat{\boldsymbol{\beta}} + \hat{\epsilon}_{[k_2]})$. We propose the following iterative algorithm to find indexes k_1 and k_2 : Let $k_1 = n_1$ and $k_2 = n_2$ be the initial values for k_1 and k_2 .

Step 1: If $\hat{g}(\hat{\epsilon}_{[k_1]}) < \hat{g}(\hat{\epsilon}_{[k_2]})$ and $\hat{g}(\hat{\epsilon}_{[k_1+1]}) < \hat{g}(\hat{\epsilon}_{[k_2+1]})$, $k_1 \leftarrow k_1 + 1$ and $k_2 \leftarrow k_2 + 1$; if $\hat{g}(\hat{\epsilon}_{[k_1]}) > \hat{g}(\hat{\epsilon}_{[k_2]})$ and $\hat{g}(\hat{\epsilon}_{[k_1-1]}) > \hat{g}(\hat{\epsilon}_{[k_2-1]})$, $k_1 \leftarrow k_1 - 1$ and $k_2 \leftarrow k_2 - 1$.

Step 2: Iterate the above procedure until none of above two conditions is satisfied or $(k_1 - 1)(k_2 - n) = 0$.

We use this method to construct prediction intervals for MODLR.

Table 3: Average widths (percentage of coverage) of the prediction intervals

Methods	Nominal confidence levels			
	10%	30%	50%	90%
LSE	2.166(0.101)	6.687(0.294)	12.70(0.493)	53.03(0.896)
MEDREG	0.975(0.091)	2.638(0.292)	6.506(0.491)	48.52(0.894)
MM-estimate	1.144(0.099)	2.910(0.294)	6.499(0.497)	48.49(0.906)
MODLRE	0.012(0.112)	0.035(0.311)	0.571(0.499)	26.44(0.899)

In Table 3, we report the average widths and the actual coverage rates of the prediction intervals for 10%, 30%, 50%, and 90% confidence levels. The actual coverage rates are estimated based on leave-one-out cross validation. From Table 3, we have the following findings:

1. All the prediction intervals are well-calibrated — the actual coverage rates are very close to the nominal confidence levels.
2. The average widths of prediction intervals constructed around the point prediction defined by MODLRE are significantly shorter than the prediction intervals constructed around the other three estimates.
3. Both MEDREG and MM-estimate have shorter prediction intervals than LSE.

4. SUMMARY AND DISCUSSIONS

In this article we proposed a new data analysis tool called modal linear regression in order to explore the relationship between a response variable and a set of predictors. Modal linear regression investigates this relationship using the conditional mode instead of the conditional mean or other summaries used by traditional regression techniques. When the error distribution is skewed, modal linear regression provides a more meaningful prediction than LSE. Our empirical results show that the modal linear regression provides significantly shorter prediction intervals than other regression methods.

In the application to the forest fire dataset, we provided one possible way to construct asymmetric prediction intervals for MODLR. Based on cross-validation results, the proposed skewed prediction intervals for MODLR were much shorter than the prediction intervals constructed by some of the other commonly used regression methods for forest fire data. Further research can be conducted to find out how to construct the shortest (skewed) prediction interval for a given confidence level using the information of skewed error density. One related work is by Kim and Lindsay (2011), who proposed to use confidence distribution sampling to visualize confidence sets.

Modal linear regression assumes that the mode of the conditional density of Y given \mathbf{x} is a linear function of \mathbf{x} . The idea of modal linear regression can be easily generalized to other models such as nonlinear regression, non-parametric regression, and varying coefficient partial linear regression. In addition, it would also be interesting to see how to select the most informative variables based on this modal regression idea. This will comprise our future research work.

ACKNOWLEDGMENTS

The authors would like to thank the Editor, an Associate Editor, and the anonymous referees for their valuable comments and suggestions, which led to substantial improvements in the

manuscript.

APPENDIX

The following technical conditions are imposed in this section.

(A1) $g^{(v)}(t | \mathbf{x})$, $v = 0, 1, 2, 3$ is continuous in a neighborhood of 0, and $g'(0 | \mathbf{x}) = 0$ for any \mathbf{x} .

(A2) $n^{-1} \sum_{i=1}^n g''(0 | \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T = J + o_p(1)$, $n^{-1} \sum_{i=1}^n g'''(0 | \mathbf{x}_i) \mathbf{x}_i = K + o_p(1)$, and $n^{-1} \sum_{i=1}^n g(0 | \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T = L + o_p(1)$, where $J < 0$, i.e., $-J$ is a positive definite matrix.

(A3) $n^{-1} \sum_{i=1}^n \|\mathbf{x}_i\|^4 = O_p(1)$.

0, and $g'(0 | \mathbf{x}) = 0$. any \mathbf{x} .

The above conditions are mild and are fulfilled in many applications. Note that the J, K , and L are defined in Theorem 2.3. All the results proved in this section also hold if general kernels are used for ϕ in (2.3) under some mild conditions adopted for traditional kernel density estimator (for example, ϕ is symmetric about 0 and has bounded continuous third derivative. In addition, ϕ has finite second moment with $\int t^2 \phi^2(t) dt < \infty$).

Proof of Theorem 2.1: Note that

$$\begin{aligned}
\log\{Q_h(\boldsymbol{\beta}^{(k+1)})\} - \log\{Q_h(\boldsymbol{\beta}^{(k)})\} &= \log \left\{ \sum_{i=1}^n \phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k+1)}) \right\} - \log \left\{ \sum_{i=1}^n \phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k)}) \right\} \\
&= \log \left[\frac{\sum_{i=1}^n \phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k+1)})}{\sum_{i=1}^n \phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k)})} \right] \\
&= \log \left[\sum_{i=1}^n \frac{\phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k)})}{\sum_{i=1}^n \phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k)})} \frac{\phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k+1)})}{\phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k)})} \right] \\
&= \log \left[\sum_{i=1}^n \pi(i | \boldsymbol{\beta}^{(k)}) \frac{\phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k+1)})}{\phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k)})} \right]
\end{aligned}$$

Based on the Jensen's inequality, we have

$$\log\{Q_h(\boldsymbol{\beta}^{(k+1)})\} - \log\{Q_h(\boldsymbol{\beta}^{(k)})\} \geq \sum_{i=1}^n \pi(i | \boldsymbol{\beta}^{(k)}) \log \left\{ \frac{\phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k+1)})}{\phi_h(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(k)})} \right\}.$$

Based on the property of M step in (2.6), we have

$$\log\{Q_h(\boldsymbol{\beta}^{(k+1)})\} - \log\{Q_h(\boldsymbol{\beta}^{(k)})\} \geq 0$$

and thus $Q_h(\boldsymbol{\beta}^{(k+1)}) \geq Q_h(\boldsymbol{\beta}^{(k)})$.

Proof of Theorem 2.2: Note that

$$\phi_h''(t) = h^{-3} \left(\frac{t^2}{h^2} - 1 \right) \phi(t/h) \text{ and } \phi_h'(t) = -\frac{t}{h^3} \phi(t/h).$$

Let $a_n = (nh^3)^{-1/2} + h^2$. It is sufficient to show that for any given $\eta > 0$, there exists a large constant c such that

$$P\left\{ \sup_{\|\mu\|=c} Q_h(\boldsymbol{\beta}_0 + a_n \mu) < Q_h(\boldsymbol{\beta}_0) \right\} \geq 1 - \eta. \quad (\text{A.1})$$

Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ and $\mathbf{y} = (y_1, \dots, y_n)^T$. Denote

$$K_n \equiv \frac{\partial Q_h(\boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta}} = -\frac{1}{n} \sum_{i=1}^n \phi_h'(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_0) \mathbf{x}_i \quad (\text{A.2})$$

$$J_n \equiv \frac{\partial^2 Q_h(\boldsymbol{\beta}_0)}{\partial \boldsymbol{\beta} \boldsymbol{\beta}^T} = \frac{1}{n} \sum_{i=1}^n \phi_h''(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_0) \mathbf{x}_i \mathbf{x}_i^T, \quad (\text{A.3})$$

where $Q_h(\boldsymbol{\beta})$ is defined in (2.3) and $\boldsymbol{\beta}_0$ is the true parameter value.

Based on Taylor expansion and symmetric property of $\phi(t)$, we can get the mean and

variance of J_n and K_n :

$$\begin{aligned}
\mathbb{E}(J_n \mid \mathbf{x}) &= \frac{1}{n} \sum_{i=1}^n g''(0 \mid \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T \{1 + o_p(1)\} = J \{1 + o_p(1)\}, \\
\text{Var}(J_n \mid \mathbf{x}) &= O_p\{(nh^5)^{-1}\}, \\
\mathbb{E}(K_n \mid \mathbf{x}) &= \frac{h^2}{2n} \sum_{i=1}^n g'''(0 \mid \mathbf{x}_i) \mathbf{x}_i (1 + o_p(1)) = \frac{h^2}{2} K \{1 + o_p(1)\}, \\
\text{Cov}(K_n \mid \mathbf{x}) &= \frac{1}{n^2 h^3} \nu_2 \sum_{i=1}^n g(0 \mid \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T \{1 + o(1)\} = \frac{1}{nh^3} \nu_2 L \{1 + o_p(1)\}, \tag{A.4}
\end{aligned}$$

where $J = \lim n^{-1} \sum_{i=1}^n g''(0 \mid \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T$, $K = \lim n^{-1} \sum_{i=1}^n g'''(0 \mid \mathbf{x}_i) \mathbf{x}_i$, and $L = \lim n^{-1} \sum_{i=1}^n g(0 \mid \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T$. By default, when calculating the variance of a matrix, we find the variance of each element of the matrix. Using the result $X = \mathbb{E}(X) + O_p(\{\text{Var}(X)\}^{1/2})$, since $nh^5 \rightarrow \infty$, $J_n = J + o_p(1)$. Notice that

$$\begin{aligned}
Q_h(\boldsymbol{\beta}_0 + a_n \boldsymbol{\mu}) - Q_h(\boldsymbol{\beta}_0) &= a_n K_n^T \boldsymbol{\mu} + \frac{a_n^2}{2} \boldsymbol{\mu}^T J_n \boldsymbol{\mu} - \frac{a_n^3}{6nh^4} \sum_{i=1}^n \phi''' \left(\frac{y_i - \mathbf{x}_i^T \boldsymbol{\beta}^*}{h} \right) (\mathbf{x}_i^T \boldsymbol{\mu})^3 \\
&= M_1 + M_2 + M_3, \tag{A.5}
\end{aligned}$$

where $\|u\| = c$ and $\|\boldsymbol{\beta}^* - \boldsymbol{\beta}_0\| \leq ca_n$. From (A.4), we get $K_n = O_p(a_n)$ and hence $M_1 = O_p(a_n^2)$. Note that $M_2 = 0.5a_n^2 \boldsymbol{\mu}^T J \boldsymbol{\mu} \{1 + o_p(1)\}$. Based on the boundness of $\phi^{(4)}(t)$ and $\|\boldsymbol{\beta}^* - \boldsymbol{\beta}_0\| \leq ca_n$, we have

$$\phi''' \left(\frac{y_i - \mathbf{x}_i^T \boldsymbol{\beta}^*}{h} \right) = \phi''' \left(\frac{y_i - \mathbf{x}_i^T \boldsymbol{\beta}_0}{h} \right) (1 + o_p(1)).$$

Noting that $\phi'''(t) = (3t - t^3)\phi(t)$, based on the Taylor expansion and the symmetric property

of $\phi(t)$, we have that

$$\mathbb{E} \left\{ \phi''' \left(\frac{Y_i - \mathbf{x}_i^T \boldsymbol{\beta}_0}{h} \right) \middle| \mathbf{x} \right\} = O_p(h^4), \quad \text{Var} \left\{ \phi''' \left(\frac{Y_i - \mathbf{x}_i^T \boldsymbol{\beta}_0}{h} \right) \middle| \mathbf{x} \right\} = O_p(h). \quad (\text{A.6})$$

Since $nh^5 \rightarrow \infty$, we can prove that $M_3 = o_p(a_n^2)$.

For any $\eta > 0$, we can choose c big enough, such that the second term M_2 dominates the other two terms in (A.5) with probability $1 - \eta$. Since $J < 0$, $Q_h(\boldsymbol{\beta}_0 + a_n \mu) - Q_h(\boldsymbol{\beta}_0) < 0$ with probability $1 - \eta$. The result of Theorem 2.2 follows. \square

Proof of Theorem 2.3: Suppose $\hat{\boldsymbol{\beta}}$ is the consistent solution to $\partial Q_h(\boldsymbol{\beta})/\partial \boldsymbol{\beta}$ found in Theorem 2.2. Based on the Taylor expansion, we have

$$0 = \frac{\partial Q_h(\hat{\boldsymbol{\beta}})}{\partial \boldsymbol{\beta}} = K_n + (J_n + L_n)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0), \quad (\text{A.7})$$

where

$$L_n = -\frac{1}{2nh^4} \sum_{i=1}^n \left[\phi''' \left(\frac{Y_i - \boldsymbol{\beta}^{*T} \mathbf{x}_i}{h} \right) \left\{ (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^T \mathbf{x}_i \right\} \mathbf{x}_i \mathbf{x}_i^T \right],$$

where $\|\boldsymbol{\beta}^* - \boldsymbol{\beta}_0\| \leq \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\|$.

Based on the result of (A.7), we have $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 = (J_n + L_n)^{-1} K_n$. Since $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\| = O_p(a_n)$, where $a_n = (nh^3)^{-1/2} + h^2$, similar to the proof of M_3 in (A.5), we have $L_n = o_p(1)$. Hence, based on (A.4), we have $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 = J^{-1} K_n (1 + o_p(1))$. Next we prove the asymptotic normality for $K_n^* = \sqrt{nh^3} K_n$.

For any unit vector $d \in \mathbb{R}^{p+1}$, we prove

$$\{d^T \text{Cov}(K_n^*) d\}^{-\frac{1}{2}} \{d^T K_n^* - d^T \mathbb{E}(K_n^*)\} \xrightarrow{D} N(0, 1)$$

Let

$$\xi_i = -\frac{1}{\sqrt{nh}}\phi'\left(\frac{Y_i - \mathbf{x}_i^T \boldsymbol{\beta}_0}{h}\right)d^T \mathbf{x}_i.$$

Then $d^T K_n^* = \sum_{i=1}^n \xi_i$. We check the Lyapunov's condition. Based on the results (A.4), we know

$$\text{Cov}(K_n) = \frac{L}{nh^3}\nu_2\{1 + o(1)\}. \quad (\text{A.8})$$

Hence $\text{Var}(d^T K_n^*) = nh^3 d^T \text{Cov}(K_n) d = g(0)\nu_2 d^T L d + o(1)$. So we only need to prove $n\text{E}|\xi_1|^3 \rightarrow 0$. Noticing that $(d^T \mathbf{x}_i)^2 \leq \|\mathbf{x}_i\|^2 \|d\|^2 = \|\mathbf{x}_i\|^2$, and $\phi'(\cdot)$ is bounded, we have $n\text{E}|\xi_1|^3 \leq O\{(nh^3)^{-1/2}\} \rightarrow 0$. So, the asymptotic normality for K_n^* holds, i.e., we have

$$\sqrt{nh^3} \{K_n - h^2 K/2(1 + o_p(h^2))\} \xrightarrow{D} N(0, \nu_2 L).$$

Based on the Slutsky's theorem, we have

$$\sqrt{nh^3} \left[\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} - \frac{h^2}{2} J^{-1} K \{1 + o_p(h^2)\} \right] \xrightarrow{D} N \{0, \nu_2 J^{-1} L J^{-1}\}. \quad \square$$

Proof of Theorem 2.4:

Since $\hat{\boldsymbol{\beta}}_s$ has root n consistency, the asymptotic result of $\hat{\boldsymbol{\beta}}_0$ is the same as if $\hat{\boldsymbol{\beta}}_s$ were known and its asymptotic distribution can be derived from Theorem 2.3 by assuming $\mathbf{x} = 1$ and the independence of ϵ and \mathbf{x} , under which we have $J^{-1}K = g'''(0)/(2g''(0))$ and $J^{-1}LJ^{-1} = g''(0)^{-2}g(0)$. Then the result follows.

Proof of Theorem 2.5: Let $\phi^*(t) = \sqrt{2\pi}h\phi_h(t)$. Then $M = \sum_{i=1}^n \phi^*(y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}})$, where $\hat{\boldsymbol{\beta}} = T(\mathbf{Z})$. Notice that $\phi^*(\cdot)$ has a maximum at 0 with $\phi^*(0) = 1$ and $\phi^*(\cdot)$ decreases monotonely toward both sides, and that $\lim \phi^*(t) = 0$ for $|t| \rightarrow \infty$.

We first prove that $T(\mathbf{Z} \cup \mathbf{Z}')$ stays bounded if $m < M$. Let $\xi > 0$ be such that

$m + n\xi < M$, and let C be such that $\phi^*(t) \leq \xi$ for $|t| \geq C$. Let $\boldsymbol{\beta}$ be any real vector such that $|y - \mathbf{x}^T \boldsymbol{\beta}| \geq C$ for all $\mathbf{z} = (\mathbf{x}, y)$ in \mathbf{Z} . Then

$$\sum_{i=1}^{m+n} \phi^*(y_i - \mathbf{x}_i^T T(\mathbf{Z})) \geq M \quad (\text{A.9})$$

and

$$\sum_{i=1}^{m+n} \phi^*(y_i - \mathbf{x}_i^T \boldsymbol{\beta}) \leq n\xi + m. \quad (\text{A.10})$$

From (A.9) and (A.10), one knows that $T(\mathbf{Z} \cup \mathbf{Z}')$ must satisfy $|y - \mathbf{x}^T T(\mathbf{Z} \cup \mathbf{Z}')| < C$ for a point in \mathbf{Z} and thus $T(\mathbf{Z} \cup \mathbf{Z}')$ is bounded.

On the other hand, if $m > M$, let $\xi > 0$ such that $m - m\xi > M$, and let C be such that $\phi^*(t) \leq \xi$ for $|t| \geq C$. Assume that all points $\{(\mathbf{x}_{n+1}, y_{n+1}), \dots, (\mathbf{x}_{n+m}, y_{n+m})\}$ in \mathbf{Z}' are the same and satisfy a linear relationship $y = \mathbf{x}^T \boldsymbol{\beta}^*$. Let $\boldsymbol{\beta}$ be any vector such that $|y_{n+1} - \mathbf{x}_{n+1}^T \boldsymbol{\beta}| < C$. Then

$$\sum_{i=1}^{m+n} \phi^*(y_i - \mathbf{x}_i^T \boldsymbol{\beta}) \leq M + m\xi, \quad (\text{A.11})$$

and

$$\sum_{i=1}^{m+n} \phi^*(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^*) \geq m. \quad (\text{A.12})$$

From (A.11) and (A.12), one knows that $T(\mathbf{Z} \cup \mathbf{Z}')$ must satisfy $|y_{n+1} - \mathbf{x}_{n+1}^T T(\mathbf{Z} \cup \mathbf{Z}')| \leq C$. If we let $y_{n+1} \rightarrow \infty$ with \mathbf{x}_{n+1} fixed, $\|T(\mathbf{Z} \cup \mathbf{Z}')\|$ must go off to infinity, and we have breakdown. \square

References

- Botev, Z. I., Grotowski, J. F., and Kroese, D. P. (2010). Kernel Density Estimation via Diffusion, *The Annals of Statistics*, 38, 2916-2957.
- Chaudhuri, P. and Marron, J. S. (1999). Sizer for Exploration of Structures in Curves. *Journal of the American Statistical Association*, 94, 807-823.
- Cortez, P. and Morais, A. (2007). A Data Mining Approach to Predict Forest Fires using Meteorological Data. *Proceedings of the 13th EPIA 2007 - Portuguese Conference on Artificial Intelligence*, December, Guimaraes, Portugal, 512-523.
- Davies, P. L. and Kovac, A. (2004). Densities, Spectral Densities and Modality. *Annals of Statistics*, 32, 1093-1136.
- Donoho, D. L. (1982). *Breakdown Properties of Multivariate Location Estimators*. Ph.D. Qualifying paper, Department of Statistics, Harvard University.
- Donoho, D. L. and Huber, P. J. (1983). The Notion of Breakdown Point. In: Doksum, K., Hodges, J. L. (Editors), *Festschrift in Honor of Erich Lehman*. Wadsworth, Belmont, CA.
- Eddy, W. F. (1980). Optimum Kernel Estimators of the Mode. *Annals of Statistics*, 8, 870-882.
- Fisher, N. I. and Marron, J. S. (2001). Mode Testing via The Excess Mass Estimate. *Biometrika*, 88, 499-517.
- Friedman J. H. and Fisher, N. I. (1999). Bump Hunting in High-Dimensional Data. *Statistics and Computing*, 9, 123-143.
- Hall, P., Minnotte, M. C., and Zhang, C. (2004). Bump Hunting with Non-Gaussian Kernels. *Annals of Statistics*, 32, 2124-2141.

- Hampel, F. R. (1971). A General Quantitative Definition of Robustness. *Annals of Mathematical Statistics*, 42, 1887-1896.
- Hampel, F. R. (1974). The Influence Curve and Its Role in Robust Estimation. *Journal of the American Statistical Association*, 69, 909-927.
- Huber, P. J. (1981). *Robust Statistics*. Wiley, New York.
- Huber, P. J. (1984). Finite Sample Breakdown Points of M- and P-Estimators. *Annals of Statistics*, 12, 119-126.
- Kemp, G. C. R. and Santos Silva, J. M. C. (2010). Regression towards the mode. Department of Economics, University of Essex, Discussion Paper No 686. <http://privatewww.essex.ac.uk/jmc/ss/research.html>.
- Kim, D. and Lindsay, B. G. (2011). Using Confidence Distribution Sampling to Visualize Confidence Sets. *Statistica Sinica*, 21(2), 923-948.
- Lee, M. J. (1989). Mode Regression *Journal of Econometrics*, 42, 337-349.
- Li, J., Ray, S., and Lindsay, B. G. (2007). A Nonparametric Statistical Approach to Clustering via Mode Identification. *Journal of Machine Learning Research*, 8(8), 1687-1723.
- Muller, D. W. and Sawitzki, G. (1991). Excess Mass Estimates and Tests for Multimodality. *Journal of the American Statistical Association*, 86, 738-746.
- Parzen, E. (1962). On Estimation of a Probability Density Function and Mode. *The Annals of Mathematical Statistics*, 33, 1065-1076.
- Ray, S. and Lindsay, B. G. (2005). The Topography of Multivariate Normal Mixtures. *Annals of Statistics*, 33, 2042-2065.

- Silverman, B. W. (1986). *Density Estimation for Statistics and Data Analysis*. Chapman and Hall: London.
- Scott, D. W. (1992). *Multivariate Density Estimation: Theory, Practice and Visualization*. New York: Wiley.
- Yao, W. (2013). A Note On EM Algorithm For Mixture Models. *Statistics and Probability Letters*, 83, 519-526.
- Yao, W. and Lindsay, B. G. (2009). Bayesian Mixture Labeling by Highest Posterior Density. *Journal of American Statistical Association*, 104, 758-767.
- Yohai, V. J. (1987). High Breakdown-point and High Efficiency Robust Estimates for Regression. *The Annals of Statistics*, 15, 642-656.