Modal Regression using Kernel Density Estimation: a Review

Yen-Chi Chen*

Abstract

We review recent advances in modal regression studies using kernel density estimation. Modal regression is an alternative approach for investigating relationship between a response variable and its covariates. Specifically, modal regression summarizes the interactions between the response variable and covariates using the conditional mode or local modes. We first describe the underlying model of modal regression and its estimators based on kernel density estimation. We then review the asymptotic properties of the estimators and strategies for choosing the smoothing bandwidth. We also discuss useful algorithms and similar alternative approaches for modal regression, and propose future direction in this field.

1 Introduction

Modal regression is an approach for studying the relationship between a response variable Y and its covariates X. Instead of seeking the conditional mean, the modal regression searches for the conditional modes (Collomb et al., 1986; Lee, 1989; Sager and Thisted, 1982) or local modes (Chen et al., 2016a; Einbeck and Tutz, 2006) of the response variable Y given the covariate X = x. The modal regression would be a more reasonable modeling approach than the usual regression in two scenarios. First, when the conditional density function is skewed or has a heavy tail. When the conditional density function has skewness, the conditional mean may not provide a good representation for summarizing the relations between the response and the covariate (X-Y) relation). The other scenario is when the conditional density function has multiple local modes. This occurs when the X-Y relation

^{*}Department of Statistics, University of Washington

contains multiple patterns. The conditional mean may not capture any of these patterns so it can be a very bad summary; see, e.g., Chen et al. (2016a) for an example. This situation has already been pointed out in Tarter and Lock (1993), where the authors argue that we should not stick to a single function for summarizing the X-Y relation and they recommend looking for the conditional local modes.

Modal regression has been applied to various problems such as predicting Alzheimer's disease (Wang et al., 2017), analyzing dietary data (Zhou and Huang, 2016), predicting temperature (Hyndman et al., 1996), analyzing electricity consumption (Chaouch et al., 2017), and studying the pattern of forest fire (Yao and Li, 2014). In particular, Wang et al. (2017) argued that the neuroimaging features and cognitive assessment are often heavy-tailed and skewed. A traditional regression approach may not work well in this scenario, so the authors propose to use a regularized modal regression for predicting Alzheimer's disease.

The concept of modal regression was proposed in Sager and Thisted (1982). In this pioneering work, the authors stipulated that the conditional (global) mode be a monotone function of the covariate. Sager and Thisted (1982) also pointed out that a modal regression estimator can be constructed using a plug-in from a density estimate. Lee (1989) proposed a linear modal regression that combined a smoothed 0-1 loss with a maximum likelihood estimator (see equation (7) for how these two ideas are connected). The idea proposed in Lee (1989) was subsequently modified in many studies; see, e.g., Kemp and Silva (2012); Krief (2017); Lee (1989, 1993); Lee and Kim (1998); Manski (1991); Yao and Li (2014).

The idea of using conditional local modes has been pointed out in Tarter and Lock (1993) and the 1992 version of Dr. David Scott's book Multivariate density estimation: theory, practice, and visualization (Scott, 1992). The first systematic analysis was done in Einbeck and Tutz (2006), where the authors proposed a plug-in estimator using a kernel density estimator (KDE) and computed their estimator by a computational approach modified from the meanshift algorithm (Cheng, 1995; Comaniciu and Meer, 2002; Fukunaga and Hostetler, 1975). The theoretical analysis and several extensions, including confidence sets, prediction sets, and regression clustering were later studied in Chen et al. (2016a). Recently, Zhou and Huang (2016) extended this idea to measurement error problems.

The remainder of this review paper is organized as follows. In Section 2, we formally

define the modal regression model and discuss its estimator by KDE. In Section 3, we review the asymptotic theory of the modal regression estimators. Possible strategies for selecting the smoothing bandwidth and computational techniques are proposed in Section 4 and 5, respectively. In Section 6, we discuss two alternative but similar approaches to modal regression – the mixture of regression and the regression quantization method. The review concludes with some possible future directions in Section 7.

2 Modal Regression

For simplicity, we assume that the covariate X is univariate with a compactly supported denisty function. Two types of modal regression have been studied in the literature. The first type, focusing on the conditional (global) mode, is called uni-modal regression (Collomb et al., 1986; Lee, 1989; Manski, 1991; Sager and Thisted, 1982). The other type, which finds the conditional local modes, is called multi-modal regression (Chen et al., 2016a; Einbeck and Tutz, 2006).

More formally, let q(z) denote the probability density function (PDF) of a random variable Z. We define the operators

$$\mathsf{UniMode}(Z) = \mathop{\mathrm{argmax}}_{z} \ q(z)$$

and

$$\mathsf{MultiMode}(Z) = \{z : q'(z) = 0, q''(z) < 0\},$$

which return the global mode and local modes of the PDF of Z, respectively. Note that we need q to be twice differentiable. Uni-modal regression searches for the function

$$m(x) = \mathsf{UniMode}(Y|X = x) = \operatorname*{argmax}_{y} \ p(y|x) \tag{1}$$

whereas multi-modal regression targets

$$M(x) = \mathsf{MultiMode}(Y|X = x) = \left\{ y : \frac{\partial}{\partial y} p(y|x) = 0, \frac{\partial^2}{\partial y^2} p(y|x) < 0 \right\}. \tag{2}$$

Note that the modal function M(x) may be a multi-valued function. Namely, M(x) may take multiple values at a given point x. Figure 1 presents examples of uni-modal and multi-modal regression using a plug-in estimate from a KDE.

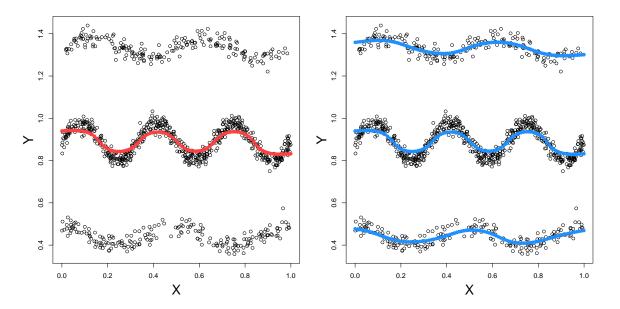


Figure 1: Uni-modal regression (left; red curve) and multi-modal regression (right; blue curves) on a simulation dataset with three components.

Because $p(y|x) = \frac{p(x,y)}{p(x)}$, the mode or local modes of p(y|x) and p(x,y) are equal for a given fixed x. Thus, provided that p(x) > 0, we can rewrite both uni-modal and multi-modal regression in the following form:

$$m(x) = \operatorname*{argmax}_{y} \ p(x,y), \quad M(x) = \left\{ y : \frac{\partial}{\partial y} p(x,y) = 0, \frac{\partial^2}{\partial y^2} p(x,y) < 0 \right\}. \tag{3}$$

That is, both types of modal regressions can be directly defined through the joint PDF. Therefore, an estimated joint PDF can be inverted into a modal regression estimate. Note that there are also Bayesian methods for modal regression; see, e.g., Ho et al. (2017).

2.1 Estimating Uni-modal Regression

The KDE provides a simple approach for estimating uni-modal regression (Collomb et al., 1986; Sager and Thisted, 1982; Yao et al., 2012). After estimating the joint PDF, we form a plug-in estimate for the uni-modal regression using KDE. In more detail, let

$$\widehat{p}_n(x,y) = \frac{1}{nh_1h_2} \sum_{i=1}^n K_1\left(\frac{X_i - x}{h_1}\right) K_2\left(\frac{Y_i - y}{h_2}\right) \tag{4}$$

be the KDE where K_1 and K_2 are kernel functions such as Gaussian functions and $h_1, h_2 > 0$ are smoothing parameters that control the amount of smoothing. An estimator of m is

$$\widehat{m}_n(x) = \underset{y}{\operatorname{argmax}} \ \widehat{p}_n(x, y). \tag{5}$$

Note that the joint PDF can be estimated by other approaches such as local polynomial estimation as well (Einbeck and Tutz, 2006; Fan et al., 1996; Fan and Yim, 2004).

Equation (4) has been generalized to the case of censored response variables. Khardani et al. (2010, 2011); Ould-Saïd and Cai (2005). Suppose that instead of observing the response variables Y_1, \dots, Y_n , we observe $T_i = \min\{Y_i, C_i\}$ and an indicator $\delta_i = I(T_i = Y_i)$ that informs whether Y_i is observed or not and C_i is an random variable that is independent of X_i and Y_i . In this case, equation (4) can be modified to

$$\widehat{p}_n^{\dagger}(x,y) = \frac{1}{nh_1h_2} \sum_{i=1}^n K_1\left(\frac{X_i - x}{h_1}\right) K_2\left(\frac{T_i - y}{h_2}\right) \times \frac{\delta_i}{\widehat{S}_n(T_i)},\tag{6}$$

where $\widehat{S}_n(t)$ is the Kaplan-Meier estimator (Kaplan and Meier, 1958)

$$\widehat{S}_n(t) = \begin{cases} \prod_{i=1}^n \left(1 - \frac{\delta_{(i)}}{n-i+1}\right)^{I(T_{(i)} \le t)} & \text{if } t < T_{(n)}, \\ 0 & \text{otherwise,} \end{cases}$$

with $T_{(1)} \leq T_{(2)} \leq \cdots \leq T_{(n)}$ being the ordered T_i 's and $\delta_{(i)}$ being the value of δ for the *i*-th ordered observation. Replacing \widehat{p}_n by \widehat{p}_n^{\dagger} in equation (5), we obtain a uni-modal regression estimator in the censoring case.

Uni-modal regression may be estimated parametrically as well. When K_2 is a spherical (box) kernel $K_2(x) = \frac{1}{2}I(|x| \le 1)$, the argmax operation is equivalent to the argmin operator on a flattened 0-1 loss. In more detail, consider a 1D toy example with observations Z_1, \dots, Z_n and a corresponding KDE $\widehat{q}(z) = \frac{1}{2nh} \sum_{i=1}^n I(|z-Z_i| \le h)$ obtained with a spherical kernel. It is easily seen that

Parametric uni-modal regression forms estimators using equation (7) or its generalizations (Kemp and Silva, 2012; Khardani and Yao, 2017; Krief, 2017; Lee, 1989, 1993; Lee and Kim, 1998; Manski, 1991; Yao and Li, 2014). Parameters estimated through the maximizing criterion in equation (7) is equivalent to maximum likelihood estimation. Conversely, parameter estimation through the minimization procedure in equation (7) is equivalent to empirical risk minimization. For example, to fit a linear model to $m(x) = \beta_0 + \beta_1 x$ (Lee, 1989; Yao and Li, 2014), we can use the fitted parameters

$$\widehat{\beta}_{0}, \widehat{\beta}_{1} = \underset{\beta_{0}, \beta_{1}}{\operatorname{argmax}} \frac{1}{2nh} \sum_{i=1}^{n} I(|\beta_{0} + \beta_{1}X_{i} - Y_{i}| \leq h)$$

$$= \underset{\beta_{0}, \beta_{1}}{\operatorname{argmin}} \sum_{i=1}^{n} I(|\beta_{0} + \beta_{1}X_{i} - Y_{i}| > h)$$

$$(8)$$

to construct our final estimate of m(x).

Using equation (7), we can always convert the problem of finding the uni-modal regression into a problem of minimizing a loss function. Here, the tuning parameter h can be interpreted as the smoothing bandwidth of the applied spherical kernel. Choosing h is a persistently difficult task. Some possible approaches will be discussed in Section 4.

2.2 Estimating Multi-modal Regression

Like uni-modal regression, multi-modal regression can be estimated using a plug-in estimate from the KDE (Chen et al., 2016a; Einbeck and Tutz, 2006). Recalling that $\hat{p}_n(x, y)$ is the KDE of the joint PDF, an estimator of M(x) is

$$\widehat{M}_n(x) = \left\{ y : \frac{\partial}{\partial y} \widehat{p}_n(x, y) = 0, \frac{\partial^2}{\partial y^2} \widehat{p}_n(x, y) < 0 \right\}.$$
(9)

Namely, we use the conditional local modes of the KDE to estimate the conditional local modes of the joint PDF. Plug-ins from a KDE have been applied in estimations of many structures (Chen, 2017; Scott, 2015) such as the regression function(Nadaraya, 1964; Watson, 1964), modes (Chacón and Duong, 2013; Chen et al., 2016b), ridges (Chen et al., 2016a; Genovese et al., 2014), and level sets (Chen et al., 2017a; Rinaldo and Wasserman, 2010). An alternative way of estimating the multi-modal regression was proposed in Sasaki et al. (2016).

In the measurement error case where the covariates X_1, \dots, X_n are observed with noises, we can replace K_1 by a deconvolution kernel to obtain a consistent estimator (Zhou and Huang, 2016). In more detail, let

$$W_i = X_i + U_i, \ i = 1, \cdots, n,$$

where U_1, \dots, U_n are IID measurement errors that are independent of the covariates and responses. We assume that the PDF of U_1 , $f_U(u)$, is known. Here we observe not X_i 's but pairs of $(W_1, Y_1), \dots, (W_n, Y_n)$. Namely, we observe the response variable and its corrupted covariate. In this case, (4) is replaced by

$$\widetilde{p}_n(x,y) = \frac{1}{nh_1h_2} \sum_{i=1}^n K_U\left(\frac{W_i - x}{h_1}\right) K_2\left(\frac{Y_i - y}{h_2}\right),$$
(10)

where

$$K_U(t) = \frac{1}{2\pi} \int e^{-its} \frac{\phi_{K_1}(s)}{\phi_U(s/h_1)} ds$$

with ϕ_{K_1} and ϕ_U being the Fourier transforms of K_1 and f_U , respectively. The estimator of M then becomes the conditional local modes of \widetilde{p} :

$$\widetilde{M}_n(x) = \left\{ y : \frac{\partial}{\partial y} \widetilde{p}_n(x, y) = 0, \frac{\partial^2}{\partial y^2} \widetilde{p}_n(x, y) < 0 \right\}$$
(11)

For more details, the reader is referred to Zhou and Huang (2016).

2.3 Uni-modal versus Multi-modal Regression

Uni-modal and multi-modal regression have their own advantages and disadvantages. Unimodal regression is an alternative approach for summarizing the covariate-response relationship using a single function. Multi-modal regression performs a similar job but allows a
multi-valued summary function. When the relation between the response and the covariate
is complicated or has several distinct components (see, e.g., Figure 1), multi-modal regression
may detect the hidden relation that cannot be found by uni-modal regression. In particular, the prediction regions tend to be smaller in multi-modal regression than in uni-modal
regression (see Figure 2 for an example). However, multi-modal regression often returns a
multi-valued function that is more difficult to interpret than the output from a uni-modal
regression.

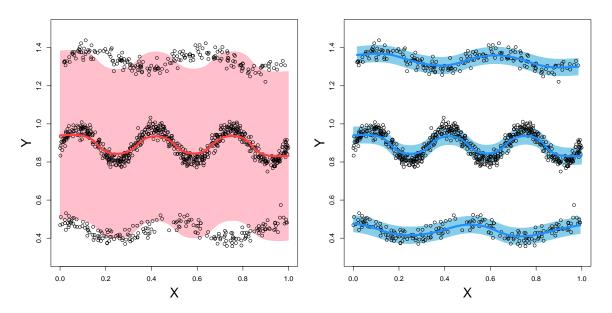


Figure 2: 90% prediction regions constructed from uni-modal regression (pink area in the left panel) and multi-modal regression (light blue area in the right panel). Clearly, the prediction region is much smaller in the multi-modal regression than in uni-modal regression because multi-modal regression detects all components whereas uni-modal regression discovers only the main component.

3 Consistency of Modal Regression

3.1 Uni-Modal Regression

Uni-modal regression often makes some smoothness assumptions¹ on the conditional density function p(y|x) over variable y (Lee, 1989). These assumptions are made to convert the mode hunting problem into a minimization or maximization problem in equation (7). Equation (7) implies that many estimators implicitly smooth the data by a spherical kernel then select the point that maximizes the result. Thus, the estimator converges to the mode of a *smoothed* density function (the expectation of the KDE). To ensure that the mode of the smoothed density function remains at the same location as the mode of the original density function, a symmetric assumption is necessary. The convergence rate of a parametric model with a box kernel for variable Y is $O_P(n^{-1/3})$ (Lee, 1989). If the box kernel is replaced by a quadratic kernel, the convergence rate becomes $O_P(1/\sqrt{n})$ under suitable assumptions (Lee, 1993). A nonparametric convergence rate was derived in Yao and Li (2014) under a weaker assumption.

When estimating m(x) using a plug-in $\widehat{m}_n(x)$ from KDE, the convergence rate depends on the assumptions. If the conditional density possesses good characteristics (such as symmetry), then

$$\widehat{m}_n(x) - m(x) = O(h_1^2) + O_P\left(\sqrt{\frac{1}{nh_1}}\right)$$
 (12)

when we are using the a Gaussian kernel or a first-order local polynomial estimator (Yao et al., 2012). Besides the convergence rate, Yao et al. (2012) also derived the asymptotic normality of the estimator:

$$\sqrt{nh_1}\left(\frac{\widehat{m}_n(x) - m(x) - h_1^2b(x)}{\sigma(x)}\right) \stackrel{D}{\to} N(0,1),$$

where b(x), $\sigma(x)$ are functions describing the asymptotic bias and variance. Note that the convergence rate and asymptotic normality are very similar to the usual nonparametric estimators. Under the assumptions on conditional density, the covariate is more responsible for the smoothing effect than the response.

¹In Lee (1989), the assumptions are either symmetric and homogeneous errors or non-symmetric and heterogeneous error.

Various studies have reported the convergence of uni-modal regression with dependent covariates (Attaoui, 2014; Collomb et al., 1986; Dabo-Niang and Laksaci, 2010; Khardani et al., 2010, 2011; Ould-Saïd, 1993, 1997; Ould-Saïd and Cai, 2005). Strong consistency was investigated in (Collomb et al., 1986; Ould-Saïd, 1993, 1997). The convergence rate has also been derived in uni-modal regression with functional dependent covariates Attaoui (2014); Dabo-Niang and Laksaci (2010), and with censored response Khardani et al. (2010, 2011); Ould-Saïd and Cai (2005).

3.2 Multi-Modal Regression

Measuring the quality of modal regression is a difficult task because the estimator \widehat{M}_n and the parameter of interest M are both multi-valued functions. $\widehat{M}_n(x)$ and M(x) are collections (sets) of values/points at each given point x.

We now define the Hausdoff distance, a popular measure of evaluating the difference between two sets. The Hausdorff distance between two sets For two given sets $A, B \subset \mathbb{R}^k$ is given by

$$\begin{aligned} \mathsf{Hausdorff}(A,B) &= \inf\{r \geq 0 : A \subset B \oplus r, B \subset A \oplus r\} \\ &= \max\left\{\sup_{x \in A} d(x,B), \sup_{x \in B} d(x,A)\right\}, \end{aligned}$$

where $A \oplus r = \{x \in \mathbb{R}^k : d(x,A) \leq r\}$ is an augmented set of A and $d(x,A) = \inf_{y \in A} ||x-y||$ is the projection distance from point x to set A. Specifically, the Hausdorff distance is the maximum projection distance between sets A and B and can be viewed as an L_{∞} distance of sets. As such, the Hausdorff distance has been applied as quality measure in estimating local modes (Chen et al., 2016b), ridges (Genovese et al., 2014), and level sets (Chen et al., 2017a), so it is excellently suitable for measuring the distance between $\widehat{M}_n(x)$ and M(x).

The pointwise error at a given point x is defined as

$$\Delta_n(x) = \mathsf{Hausdorff}\left(\widehat{M}_n(x), M(x)\right).$$

This pointwise error is similar to the usual pointwise error of estimating a regression function. Based on the pointwise error, we can easily define the *mean integrated square error (MISE)* and uniform errors

$$\mathsf{MISE}_n = \int \Delta_n^2(x) dx, \qquad \Delta_n = \sup_x \Delta_n(x).$$

These quantities are generalized from the errors in nonparametric literature (Scott, 2015).

The convergence rate of $\widehat{M}_n(x)$ has been derived in Chen et al. (2016a):

$$\Delta_{n}(x) = O(h_{1}^{2} + h_{2}^{2}) + O_{P}\left(\sqrt{\frac{1}{nh_{1}h_{2}^{3}}}\right)$$

$$\Delta_{n} = O(h_{1}^{2} + h_{2}^{2}) + O_{P}\left(\sqrt{\frac{\log n}{nh_{1}h_{2}^{3}}}\right)$$

$$MISE_{n} = O(h_{1}^{4} + h_{2}^{4}) + O\left(\frac{\log n}{nh_{1}h_{2}^{3}}\right).$$
(13)

The bias is now contributed by smoothing covariates and response. The convergence rate of the stochastic variation, $O_P\left(\sqrt{\frac{1}{nh_1h_2^3}}\right)$, depends on the amount of smoothing in both the covariate and response variable as well. The component h_2^3 can be decomposed as $h_2 \cdot h_2^2$, where the first part h_2 is the usual smoothing and the second part, h_2^2 , is from derivative estimations. Note that the convergence rates in equation (13) require no symmetric-like assumption on the conditional density, but only smoothness and bounded curvature at each local mode. Therefore, the assumptions ensuring a consistent estimator are much weaker in multi-modal regression than in uni-modal regression. However, the convergence rate is much slower in multi-modal regression than in uni-modal regression (equation (12)). Note that under the same weak assumptions as multi-modal regression, uni-modal regression can also be consistently estimated by a KDE and the convergence rate will be the same as equation (13).

Chen et al. (2016a) also derived the asymptotic distribution and a bootstrap theory of Δ_n . When we ignore the bias, the uniform error converges to the maximum of a Gaussian process and the distribution of this a maximum can be approximated by the empirical bootstrap (Efron, 1979). Therefore, by applying the bootstrap, one can construct a confidence band for the modal regression.

In the case of measurement errors, a similar convergence rate to equation (13) can also be derived under suitable conditions (Zhou and Huang, 2016). Note that in this case, the distribution of measurement errors also affects the estimation quality.

4 Bandwidth Selection

Modal regression estimation often involves some tuning parameters. In a parametric model, we have to choose a window size h in equation (7). In other models, we often require two smoothing bandwidths: one for the response variable, the other for the covariate. Here we briefly summarize some bandwidth selectors proposed in the literature.

4.1 Plug-in Estimate

In Yao et al. (2012), uni-modal regression was estimated by a local polynomial estimator. One advantage of uni-modal regression is the closed-form expression of the first-order error. Therefore, we can use a plug-in approach to obtain an initial error estimate and convert it into a possible smoothing bandwidth. This approach is very similar to the plug-in bandwidth selection in the density estimation problem (Sheather, 2004).

This approach was designed to optimally estimating the error of a uni-modal regression estimate. However, this method is often not applicable to multi-modal regression because a closed-form expression of the first-order error is often unavailable. Moreover, this approach requires a pilot estimate for the first error. If the pilot estimate is unreliable, the performance of the bandwidth selector may be seriously compromised.

4.2 Adapting from Conditional Density Estimation

A common approach for selecting the tuning parameter is based on optimizing the estimation accuracy of conditional density function (Fan et al., 1996; Fan and Yim, 2004). For instance, the authors of Einbeck and Tutz (2006) adapted the smoothing bandwidth to multi-modal regression by optimizing the conditional density estimation rate and the Silverman's normal reference rule (Silverman, 1986).

The principle of adapting from estimating the conditional density often relies on opti-

mizing the integrated squared-errors:

$$ISE = \int \int (\widehat{p}_n(y|x) - p(y|x))^2 p(x)\omega(x)dxdy$$

$$= \int \int \widehat{p}_n^2(y|x)p(x)\omega(x)dxdy$$

$$-2\int \int \widehat{p}_n(y|x)p(y|x)p(x)\omega(x)dxdy + \int \int p^2(y|x)p(x)\omega(x)dxdy,$$

where $\omega(x)$ is a user-selected weight function. For simplicity, one can choose $\omega(x) = 1$ over the range of interest. Note that in density estimation literature, bandwidth selection by this expansion is called the CV criterion (Sheather, 2004).

However, the ISE involves unknown quantities so it must be estimated. Depending on the estimating procedure, there are many other approaches such as the regression-based approach, bootstrap method, and cross-validation approach; see Zhou and Huang (2017) for a comprehensive review.

Although this approach is simple and elegant, a good density estimator does not guarantee a good estimator of the local modes. As is seen in equation (13), the convergence rate is actually slower when estimating local modes than when estimating density.

4.2.1 CV-SIMEX method

The bandwidth selection method in Zhou and Huang (2016), designed for measurement errors, combines density estimation CV with simulation extrapolation (SIMEX; Cook and Stefanski 1994).

Because the last quantity in the ISE is independent of the tuning parameter, and p(x)dx = dF(x) and p(y|x)p(x)dxdy = dF(x,y) can be replaced by their empirical versions $d\widehat{F}_n(x)$ and $d\widehat{F}_n(x,y)$, the CV criterion can be estimated by

$$CV(h_1, h_2) = \frac{1}{n} \sum_{i=1}^{n} \int \widehat{p}_{-i,n}^2(y|X_i)\omega(X_i)dy - \frac{2}{n} \sum_{i=1}^{n} \omega(X_i)\widehat{p}_{-i,n}(Y_i|X_i), \tag{14}$$

where $\widehat{p}_{-i,n}(y|x)$ is the estimated conditional density without *i*-th observation (leave *i*-th observation out). If the covariates X_1, \dots, X_n are known, we can choose h_1 and h_2 by minimizing equation (14).

Measurement errors manifest as noise in the corrupted covariates W_1, \dots, W_n . In the CV-SIMEX approach (Zhou and Huang, 2016), h_2 is determined by Siverman's rule (Silverman,

1986) and we give a brief summary for the selection of h_1 as follows. We first generate W_1^*, \dots, W_n^* where $W_i^* = W_i + U_i^*$ and U_1^*, \dots, U_n^* are IID from the measurement error distribution. Then we construct the estimator \widehat{p}_n^* by replacing X_1, \dots, X_n by W_1^*, \dots, W_n^* . We modify equation (14) by replacing X_1, \dots, X_n by W_1, \dots, W_n and replacing \widehat{p}_n by \widehat{p}_n^* Note that the weight ω will also be updated according to the range of W's. Let $CV^*(h_1)$ be the resulting CV criterion. Now we compute another CV criterion as follows. We generate $W_1^{**}, \dots, W_n^{**}$ where $W_i^{**} = W_i^* + U_i^{**}$ and $U_1^{**}, \dots, U_n^{**}$ are IID from the measurement error distribution. Similar to the previous steps, we compute a new (conditional) density estimator \widehat{p}_n^{**} by replacing X_1, \dots, X_n by $W_1^{**}, \dots, W_n^{**}$. To obtain a new CV criterion, we again modify equation (14) by replacing X_1, \dots, X_n by W_1^*, \dots, W_n^* and replacing \widehat{p}_n by \widehat{p}_n^{**} . This leads to a new CV criterion which we denoted as $CV^{**}(h_1)$. We then repeat the above process multiple times and calculate the average $\overline{CV}(h_1)$ and $\overline{CV}^{**}(h_1)$. Then we choose h_1^* to be the minimizer of $\overline{CV}^*(h_1)$ and h_1^{**} to be the minimizer of $\overline{CV}^{**}(h_1)$. The final choice of smoothing bandwidth is $\widehat{h}_1 = \frac{h_1^{**}}{h_1^{**}}$.

The CV procedure assesses the quality of estimating the conditional density. The simulation process (SIMEX part) exploits the similarity between the optimal h_1 - h_1^* relation and the h_1^* - h_1^{**} relation, i.e., $\frac{h_1,opt}{h_1^*} \approx \frac{h_1^*}{h_1^{**}}$. Thus, the smoothing bandwidth is selected by equating this approximation.

However, CV-SIMEX optimizes the quality of estimating the conditional density, not the conditional local modes. As confirmed in equation (13), the optimal convergence rate differs between density and local modes estimation, so this choice would undersmooth the local modes estimation.

4.2.2 Modal CV-criteria

Zhou and Huang (2017) proposed a generalization of the density estimation CV criterion to the multi-modal regression. The idea is to replace the ISE by

$$ISE_M = \int \mathsf{Hausdorff}^2\left(\widehat{M}_n(x), M(x)\right) p(x)\omega(x)dx$$

and find an estimator of the above quantity. Optimizing the corresponding estimated ISE_M leads to a good rule for selecting the smoothing bandwidth. In particular, Zhou and Huang

(2017) proposed to use a bootstrap approach to estimate ISE_M . The quantity ISE_M is directly tailored to the modal regression rather than conditional density estimation so it reflects the actual accuracy of modal regression.

4.3 Prediction Band Approach

Another bandwidth selector for multi-modal regression was proposed in Chen et al. (2016a). This approach optimizes the size of the prediction bands using a cross-validation (CV) in the regression analysis. After selecting a prediction level (e.g., 95%), the data is split into a training set and a validation set. The modal regression estimator is constructed from the data in the training set, and the residuals of the observations are derived from the validation set. The residual of a pair X_{val}, Y_{val} is based on the shortest distance to the nearly conditional local mode. Namely, the residual for an observation (X_{val}, Y_{val}) in the validation set is $e_{val} = \min_{y \in \widehat{M}_n(X_{val})} ||Y_{val} - y||$. The 95% quantile of the residuals specifies the radius of a 95% prediction band. The width of a prediction band is twice the estimated radius. After repeating this procedure several times as the usual CV, we obtain an average size (volume of the prediction band) of the prediction band of each smoothing bandwidth. The smoothing bandwidth is chosen to be the one that has the smallest (in terms of volume) prediction band.

When h is excessively small, there will be many conditional local modes which leads to a large prediction band. On the other hand, when h is too large, the number of conditional local modes are small but each of them has a very large band, yielding a large total size of the prediction band. Consequently, this approach leads to a stable result.

However, the prediction band approach is beset with several problems. First, the optimal choice of smoothing bandwidth depends on the prediction level, which cannot be definitely selected at present. Second, calculating the size of a band is computationally challenging in high dimensions. Third, there is no theoretical guarantee that the selected bandwidth follows the optimal convergence rate.

Note that Zhou and Huang (2017) proposed a modified CV criterion for approximating the size of prediction band without specifying the prediction level. This approach avoids the problem of selecting a prediction level and is computationally more feasible.

5 Computational Methods

In the modal regression estimation, a closed-form solution to the estimator is often unavailable. Therefore, the estimator must be computed by a numerical approach. The parameters in uni-modal regression with a parametric model can be estimated by a mode-hunting procedure (Lee, 1989). When estimating uni-modal regression by a nonparametric approach, the conditional mode can be found by the gradient ascent method or an EM-algorithm (Yao and Li, 2014; Yao et al., 2012).

The conditional local modes in multi-modal regression can also be found by a gradient ascent approach. If the kerne function of the response variable K_2 has a nice form such as being a Gaussian function, gradient ascent can be easily performed by a simple algorithm called meanshift algorithm (Cheng, 1995; Comaniciu and Meer, 2002; Fukunaga and Hostetler, 1975). In the following, we briefly review the EM algorithm and the meanshift algorithm for finding modes and explain their applications to modal regression.

5.1 EM Algorithm

A common approach for finding uni-modal regression is the EM algorithm (Dempster et al., 1977; Wu, 1983). In the case of modal regression, we use the idea from a modified method called the modal EM algorithm (Li et al., 2007; Yao et al., 2012). For simplicity, we illustrate the EM algorithm using the linear uni-modal regression problem with a single covariate (Yao and Li, 2014). Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be the observed data and recall that the uni-modal regression finds the parameters using

$$\widehat{\beta}_0, \widehat{\beta}_1 = \underset{\beta_0, \beta_1}{\operatorname{argmax}} \ \frac{1}{nh} \sum_{i=1}^n K\left(\frac{Y_i - \beta_0 - \beta_1 X_i}{h}\right). \tag{15}$$

Note that when we take $K(x) = K_2(x) = \frac{1}{2}I(|x| \le 1)$, we obtain equation (8).

Given an initial choice of parameters $\beta_0^{(0)}, \beta_1^{(0)}$, the EM algorithm iterates the following two steps until convergence $(t = 1, 2, \cdots)$:

• **E-step.** Given $\beta_0^{(t-1)}, \beta_1^{(t-1)}$, compute the weights

$$\pi\left(i|\beta_0^{(t-1)}, \beta_1^{(t-1)}\right) = \frac{K\left(\frac{Y_i - \beta_0^{(t-1)} - \beta_1^{(t-1)} X_i}{h}\right)}{\sum_{j=1}^n K\left(\frac{Y_j - \beta_0^{(t-1)} - \beta_1^{(t-1)} X_j}{h}\right)}$$

for each $i = 1, \dots, n$.

• M-step. Given the weights, update the parameters by

$$\widehat{\beta}_0^{(t)}, \widehat{\beta}_1^{(t)} = \underset{\beta_0, \beta_1}{\operatorname{argmax}} \ \frac{1}{n} \sum_{i=1}^n \pi(i | \beta_0^{(t-1)}, \beta_1^{(t-1)}) \log K\left(\frac{Y_i - \beta_0 - \beta_1 X_i}{h}\right).$$

When the kernel function K is a Gaussian, the M-step has a closed-form expression:

$$\widehat{\beta}_0^{(t)}, \widehat{\beta}_1^{(t)} = (\mathbb{X}^T \mathbb{W}_{(t)} \mathbb{X}^T)^{-1} \mathbb{X}^T \mathbb{W}_{(t)} \mathbb{Y},$$

where $\mathbb{X}^T = ((1, X_1)^T, (1, X_2)^T, \dots, (1, X_n)^T)$ is the transpose of the covariate matrix in regression problem and $\mathbb{W}_{(t)}$ is an $n \times n$ diagonal matrix with elements

$$\pi(1|\beta_0^{(t-1)},\beta_1^{(t-1)}),\cdots,\pi(n|\beta_0^{(t-1)},\beta_1^{(t-1)})$$

and $\mathbb{Y} = (Y_1, \dots, Y_n)^T$ is the response vector. This is because the problem reduces to a weighted least square estimator in linear regression. Thus, the updates can be done very quickly.

Note that the EM algorithm may stuck at the local optima (Yao and Li, 2014) so the choice of initial parameters is very important. In practice, we would recommend to rerun the EM algorithm with many different initial parameters to avoid the problem of falling in a local maximum.

The EM algorithm can be extended to nonparametric uni-modal regression as well. See Yao et al. (2012) for an example of applying the EM algorithm to find the uni-modal regression using a local polynomial estimator.

5.2 Meanshift Algorithm

To illustrate the principle of the meanshift algorithm (Cheng, 1995; Comaniciu and Meer, 2002; Fukunaga and Hostetler, 1975), we return to the 1D toy example. Suppose that we

observe IID random samples $Z_1, \dots, Z_n \sim q$. Let \widehat{q} be a KDE with a Gaussian kernel $K_G(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$. A powerful feature of the Gaussian kernel is that its nicely behave derivative:

$$K'_G(x) = -x \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2} = -x \cdot K_G(x).$$

The derivative of the KDE is then

$$\widehat{q}'(z) = \frac{d}{dz} \frac{1}{nh} \sum_{i=1}^{n} K_G \left(\frac{Z_i - z}{h} \right)
= \frac{1}{nh^3} \sum_{i=1}^{n} (Z_i - z) K_G \left(\frac{Z_i - z}{h} \right)
= \frac{1}{nh^3} \sum_{i=1}^{n} Z_i K_G \left(\frac{Z_i - z}{h} \right) - \frac{z}{nh^3} \sum_{i=1}^{n} K_G \left(\frac{Z_i - z}{h} \right).$$

Multiplying both sides by nh^3 and dividing them by $\sum_{i=1}^n K_G\left(\frac{Z_i-z}{h}\right)$, the above equation becomes

$$\frac{nh^3}{\sum_{i=1}^n K_G\left(\frac{Z_i-z}{h}\right)} \cdot \widehat{q}'(z) = \frac{\sum_{i=1}^n Z_i K_G\left(\frac{Z_i-z}{h}\right)}{\sum_{i=1}^n K_G\left(\frac{Z_i-z}{h}\right)} - z.$$

Rearranging this expression, we obtain

$$\underbrace{z}_{\text{current location}} + \underbrace{\frac{nh^3}{\sum_{i=1}^n K_G\left(\frac{Z_i-z}{h}\right)} \cdot \widehat{q}'(z)}_{\text{gradient aescent}} = \underbrace{\frac{\sum_{i=1}^n Z_i K_G\left(\frac{Z_i-z}{h}\right)}{\sum_{i=1}^n K_G\left(\frac{Z_i-z}{h}\right)}}_{\text{next location}}$$

Namely, given a point z, the value of $\frac{\sum_{i=1}^{n} Z_{i}K_{G}\left(\frac{Z_{i}-z}{h}\right)}{\sum_{i=1}^{n} K_{G}\left(\frac{Z_{i}-z}{h}\right)}$ is a shifted location by applying a gradient ascent with amount $\frac{nh^{3}}{\sum_{i=1}^{n} K_{G}\left(\frac{Z_{i}-z}{h}\right)} \cdot \widehat{q}'(z)$. Therefore, the meanshift algorithm updates an initial point $z^{(t)}$ as

$$z^{(t+1)} = \frac{\sum_{i=1}^{n} Z_i K_G \left(\frac{Z_i - z^{(t)}}{h}\right)}{\sum_{i=1}^{n} K_G \left(\frac{Z_i - z^{(t)}}{h}\right)}$$

for $t = 0, 1, \cdots$. According to the above derivation, this update moves points by a gradient ascent. Thus, the stationary point $z^{(\infty)}$ will one of the local modes of the KDE. Note that although some initial points do not converge to a local modes, these points forms a set with 0 Lebesgue measure, so can be ignored (Chen et al., 2017b).

To generalize the meanshift algorithm to multi-modal regression, we fix the covariate and shift only the response variable. More specifically, given a pair of point $(x, y^{(0)})$, we fix the

covariate value x and update the response variable as follows:

$$y^{(t+1)} = \frac{\sum_{i=1}^{n} Y_i K_1 \left(\frac{X_i - x}{h_1}\right) K_2 \left(\frac{Y_i - y^{(t)}}{h_2}\right)}{\sum_{i=1}^{n} K_1 \left(\frac{X_i - x}{h_1}\right) K_2 \left(\frac{Y_i - y^{(t)}}{h_2}\right)},$$
(16)

for $t=0,1,\cdots$. Here $K_2=K_G$ is the Gaussian kernel although the meanshift algorithm accommodates other kernel functions; see Comaniciu and Meer (2002) for a discussion. The update in equation (16) is called the conditional meanshift algorithm in Einbeck and Tutz (2006) and the partial meanshift algorithm in Chen et al. (2016a). The conditional local modes include the stationary points $y^{(\infty)}$. To find all conditional local modes, we often start with multiple initial locations of the response variable and apply equation (16) to each of them.

The kernel function for the covariate K_1 in equation (16) is not limited to a Gaussian kernel. K_1 can even be a deconvolution kernel in the presence of measurement errors (Zhou and Huang, 2016).

5.3 Available Softwares

There are many statistical packages in R for modal regression. On CRAN, there are two packages that contain functions for modal regression:

- hdrcde: the function modalreg computes a multi-modal regression using the method of Einbeck and Tutz (2006).
- 1pme: the function modereg performs a multi-modal regression that can be used in situations with or without measurement errors. This package is based on the methods in Zhou and Huang (2016). Moreover, it also has two functions for bandwidth selection moderegbw and moderegbwSIMEX that apply the bandwidth selectors described in Zhou and Huang (2017) and Zhou and Huang (2016).

Note that there is also an R package on github for modal regression: https://github.com/yenchic/ModalRegression. This package is based on the method of Chen et al. (2016a).

6 Similar Approaches to Modal Regression

Multi-modal regression is a powerful tool for detecting multiple components of the conditional density. The right panel of Figure 2 demonstrates the power of a compact prediction set obtained by multi-modal regression. Multiple components of the conditional density can also be obtained by other approaches such as the mixture of regression and regression quantization method. These approaches are briefly described below.

6.1 Mixture of Regression

When multiple components reside in the conditional density, traditional regression analysis applies a mixture of regression model (Lindsay, 1995; Quandt, 1972; Quandt and Ramsey, 1978). A general form of a mixture of regression model (Huang et al., 2013; Huang and Yao, 2012) is

$$Y|X = x \sim \sum_{\ell=1}^{L} \pi_{\ell}(x) N(m_{\ell}(x), \sigma_{\ell}^{2}(x)),$$

where $\pi_{\ell}(x) \geq 0$ is the proportion of the ℓ -th component (note that $\sum_{\ell=1}^{L} \pi_{\ell}(x) = 1$) and $m_{\ell}(x), \sigma_{\ell}^{2}(x)$ denote the mean and variance. Here we assume that the data comprises L mixtures of Gaussian components (note that this assumption can be relaxed as well). In this case, the parameter functions $\pi_{\ell}(x)$, $m_{\ell}(x)$, and $\sigma_{\ell}^{2}(x)$ are parameters of interest and must be estimated from the data. The parameter functions can be estimated using smoothing techniques and maximum likelihood estimation Huang et al. (2013).

Although the mixture of regression approach is flexible, it has several limitations. First is the identifiability problem; different combinations of the parameter functions may lead to the same or similar conditional density, which destabilizes the estimator. Second, the number of components L, must be known a priori. If we assume a parametric model for the parameter functions, L can be chosen by a model selection criterion such as the Akaike or Bayesian information criterion (Huang et al., 2013). However, assuming a parametric form decreases the flexibility of the model. Moreover, computing the estimators of parameter functions often requires an EM-algorithm, which may need several re-initializations of the initial condition to get a desired estimate.

6.2 Regression Quantization

Alternatively, the authors of Loubes and Pelletier (2017) detected multiple components in a conditional density function by combining k-means (vector quantization; Gersho and Gray 2012; Graf and Luschgy 2007) algorithm and k-nearest neighbor (kNN) approach. Their method is called regression quantization. To illustrate the idea, we consider a 1D Gaussian mixture model with L distinct components. If the components are well-separated and their proportions are similar, the a k-means algorithm with k in k will return k points (called centers in the k-means literature) that approximate the centers of Gaussians. Thus, the centers of k-means correspond to the centers of components in our data.

In a regression setting, the k-means algorithm is combined with kNN. To avoid conflict in the notations, we denote the number of centers in the k-means by L, although the algorithm itself is called k-means. For a given point x, we find those X_i 's within the k-nearest neighborhood of x, process their corresponding responses by the k-means algorithm. Let

$$W_{n,i}(x) = \begin{cases} \frac{1}{k} & \text{if } X_i \text{ is among the } k \text{ nearest neighbor of } x \\ 0 & \text{otherwise} \end{cases}$$

be the weight of each observation. Given a point x, the estimator in Loubes and Pelletier (2017) was defined as

$$\widehat{c}_1(x), \cdots, \widehat{c}_L(x) = \operatorname*{argmin}_{c_1, \cdots, c_L} \sum_{i=1}^n \min_{j=1, \cdots, L} W_{n,i}(x) \|Y_i - c_j\|^2.$$

Namely, we apply the k-means algorithm to the response variable of the k-NN observations. For correct choices of k and L, the resulting estimators properly summarize the data. Because k behaves like the smoothing parameter in the KDE, the choice of $k = k_n$ has been theoretically analyzed (Loubes and Pelletier, 2017). However, the choice of L often relies on prior knowledge about the data (number of components), although L can be chosen by a gap heuristic approach (Tibshirani et al., 2001).

7 Discussion

This paper reviewed common methods for fitting modal regressions. We discussed both uni-modal and multi-modal approaches, along with relevant topics such as large sample

theories, bandwidth selectors, and computational recipes. Here we outline some possible future directions of modal regression.

- Multi-modal regression in complex processes. Although the behavior of unimodal regression has been analyzed in dependent and censored scenarios (Collomb et al., 1986; Khardani et al., 2010, 2011; Ould-Saïd and Cai, 2005), the behavior of multi-modal regression is still unclear. The behavior of multi-modal regression in censoring cases remains an open question. Moreover, by comparing the estimator in censored response variable cases and measurement error cases (equation (6) and (10), respectively), we find that to account for censoring, we must change the kernel function of the response variable, whereas to adjust the measurement errors, we need to modify the kernel function of the covariate. Thus, we believe that the KDE can be modified to solve the censoring and measurement error problems at the same time.
- Valid confidence band. In Chen et al. (2016a), the confidence band of the multimodal regression was constructed by a bootstrap approach. However, because this confidence band does not correct bias in KDE, it requires an undersmoothing assumption. Recently, Calonico et al. (2017) proposed a debiased approach that constructs a bootstrap nonparametric confidence set without undersmoothing. The application of this approach to modal regression is another possible future direction.
- Conditional bump hunting. A classical problem in nonparametric statistics is bump hunting (Burman and Polonik, 2009; Good and Gaskins, 1980; Hall et al., 2004), which detects the number of significant local modes. In modal regression analysis, the bump hunting problem may be studied in a regression setting. More specifically, we want to detect the number of significant local modes of the conditional density function. The output will be an integer function of the covariate that informs how the number of significant local modes changes over different covariate values.

Acknowledgement

We thank two referees and the editor for their very helpful comments. Yen-Chi Chen is supported by NIH grant number U01 AG016976.

References

- Attaoui, S. (2014). On the nonparametric conditional density and mode estimates in the single functional index model with strongly mixing data. Sankhya A, 76(2):356–378.
- Burman, P. and Polonik, W. (2009). Multivariate mode hunting: Data analytic tools with measures of significance. *Journal of Multivariate Analysis*, 100(6):1198–1218.
- Calonico, S., Cattaneo, M. D., and Farrell, M. H. (2017). On the effect of bias estimation on coverage accuracy in nonparametric inference. *Journal of the American Statistical Association*, (just-accepted).
- Chacón, J. and Duong, T. (2013). Data-driven density derivative estimation, with applications to nonparametric clustering and bump hunting. *Electronic Journal of Statistics*, 7:499–532.
- Chaouch, M., Laïb, N., and Louani, D. (2017). Rate of uniform consistency for a class of mode regression on functional stationary ergodic data. *Statistical Methods & Applications*, 26(1):19–47.
- Chen, Y.-C. (2017). A tutorial on kernel density estimation and recent advances. *Biostatistics & Epidemiology*, 1(1):161–187.
- Chen, Y.-C., Genovese, C. R., Tibshirani, R. J., and Wasserman, L. (2016a). Nonparametric modal regression. *The Annals of Statistics*, 44(2):489–514.
- Chen, Y.-C., Genovese, C. R., and Wasserman, L. (2016b). A comprehensive approach to mode clustering. *Electronic Journal of Statistics*, 10(1):210–241.
- Chen, Y.-C., Genovese, C. R., and Wasserman, L. (2017a). Density level sets: Asymptotics, inference, and visualization. *Journal of the American Statistical Association*, pages 1–13.
- Chen, Y.-C., Genovese, C. R., and Wasserman, L. (2017b). Statistical inference using the morse-smale complex. *Electronic Journal of Statistics*, 11(1):1390–1433.

- Cheng, Y. (1995). Mean shift, mode seeking, and clustering. *Pattern Analysis and Machine Intelligence*, *IEEE Transactions on*, 17(8):790–799.
- Collomb, G., Härdle, W., and Hassani, S. (1986). A note on prediction via estimation of the conditional mode function. *Journal of Statistical Planning and Inference*, 15:227–236.
- Comaniciu, D. and Meer, P. (2002). Mean shift: A robust approach toward feature space analysis. Pattern Analysis and Machine Intelligence, IEEE Transactions on, 24(5):603–619.
- Cook, J. R. and Stefanski, L. A. (1994). Simulation-extrapolation estimation in parametric measurement error models. *Journal of the American Statistical association*, 89(428):1314–1328.
- Dabo-Niang, S. and Laksaci, A. (2010). Note on conditional mode estimation for functional dependent data. *Statistica*, 70(1):83–94.
- Dempster, A., Laird, N., and Rubin, D. (1977). Maximum likelihood from incomplete data via the sems algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38.
- Efron, B. (1979). Bootstrap methods: Another look at the jackknife. *Annals of Statistics*, 7(1):1–26.
- Einbeck, J. and Tutz, G. (2006). Modelling beyond regression functions: an application of multimodal regression to speed–flow data. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 55(4):461–475.
- Fan, J., Yao, Q., and Tong, H. (1996). Estimation of conditional densities and sensitivity measures in nonlinear dynamical systems. *Biometrika*, 83(1):189–206.
- Fan, J. and Yim, T. H. (2004). A crossvalidation method for estimating conditional densities. Biometrika, 91(4):819–834.

- Fukunaga, K. and Hostetler, L. (1975). The estimation of the gradient of a density function, with applications in pattern recognition. *Information Theory, IEEE Transactions on*, 21(1):32–40.
- Genovese, C. R., Perone-Pacifico, M., Verdinelli, I., and Wasserman, L. (2014). Nonparametric ridge estimation. *The Annals of Statistics*, 42(4):1511–1545.
- Gersho, A. and Gray, R. M. (2012). *Vector quantization and signal compression*, volume 159. Springer Science & Business Media, Berlin/Heidelberg, Germany.
- Good, I. and Gaskins, R. (1980). Density estimation and bump-hunting by the penalized likelihood method exemplified by scattering and meteorite data. *Journal of the American Statistical Association*, 75(369):42–56.
- Graf, S. and Luschgy, H. (2007). Foundations of quantization for probability distributions. Springer, New York, NY.
- Hall, P., Minnotte, M. C., and Zhang, C. (2004). Bump hunting with non-gaussian kernels. The Annals of Statistics, 32(5):2124–2141.
- Ho, C.-s., Damien, P., and Walker, S. (2017). Bayesian mode regression using mixtures of triangular densities. *Journal of Econometrics*, 197(2):273–283.
- Huang, M., Li, R., and Wang, S. (2013). Nonparametric mixture of regression models.

 Journal of the American Statistical Association, 108(503):929–941.
- Huang, M. and Yao, W. (2012). Mixture of regression models with varying mixing proportions: A semiparametric approach. *Journal of the American Statistical Association*, 107(498):711–724.
- Hyndman, R. J., Bashtannyk, D. M., and Grunwald, G. K. (1996). Estimating and visualizing conditional densities. *Journal of Computational and Graphical Statistics*, 5(4):315–336.
- Kaplan, E. L. and Meier, P. (1958). Nonparametric estimation from incomplete observations.

 Journal of the American statistical association, 53(282):457–481.

- Kemp, G. C. and Silva, J. S. (2012). Regression towards the mode. *Journal of Econometrics*, 170(1):92–101.
- Khardani, S., Lemdani, M., and Saïd, E. O. (2010). Some asymptotic properties for a smooth kernel estimator of the conditional mode under random censorship. *Journal of the Korean Statistical Society*, 39(4):455–469.
- Khardani, S., Lemdani, M., and Saïd, E. O. (2011). Uniform rate of strong consistency for a smooth kernel estimator of the conditional mode for censored time series. *Journal of Statistical Planning and Inference*, 141(11):3426–3436.
- Khardani, S. and Yao, A. F. (2017). Non linear parametric mode regression. *Communications in Statistics-Theory and Methods*, 46(6):3006–3024.
- Krief, J. M. (2017). Semi-linear mode regression. The Econometrics Journal, 20(2):149–167.
- Lee, M.-j. (1989). Mode regression. Journal of Econometrics, 42(3):337–349.
- Lee, M.-J. (1993). Quadratic mode regression. Journal of Econometrics, 57(1-3):1-19.
- Lee, M.-J. and Kim, H. (1998). Semiparametric econometric estimators for a truncated regression model: A review with an extension. *Statistica Neerlandica*, 52(2):200–225.
- 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.
- Lindsay, B. G. (1995). Mixture models: theory, geometry and applications. In NSF-CBMS regional conference series in probability and statistics, pages i–163. JSTOR.
- Loubes, J.-M. and Pelletier, B. (2017). Prediction by quantization of a conditional distribution. *Electronic Journal of Statistics*, 11(1):2679–2706.
- Manski, C. (1991). Regression. Journal of Economic Literature, 29(1):34–50.
- Nadaraya, E. A. (1964). On estimating regression. Theory of Probability & Its Applications, 9(1):141–142.

- Ould-Saïd, E. (1993). Estimation non paramétrique du mode conditionnel. application à la prévision. Comptes rendus de l'Académie des sciences. Série 1, Mathématique, 316(9):943–947.
- Ould-Saïd, E. (1997). A note on ergodic processes prediction via estimation of the conditional mode function. *Scandinavian journal of statistics*, 24(2):231–239.
- Ould-Saïd, E. and Cai, Z. (2005). Strong uniform consistency of nonparametric estimation of the censored conditional mode function. *Nonparametric Statistics*, 17(7):797–806.
- Quandt, R. E. (1972). A new approach to estimating switching regressions. *Journal of the American statistical association*, 67(338):306–310.
- Quandt, R. E. and Ramsey, J. B. (1978). Estimating mixtures of normal distributions and switching regressions. *Journal of the American statistical Association*, 73(364):730–738.
- Rinaldo, A. and Wasserman, L. (2010). Generalized density clustering. *The Annals of Statistics*, 38(5):2678–2722.
- Sager, T. W. and Thisted, R. A. (1982). Maximum likelihood estimation of isotonic modal regression. *The Annals of Statistics*, 10(3):690–707.
- Sasaki, H., Ono, Y., and Sugiyama, M. (2016). Modal regression via direct log-density derivative estimation. In *International Conference on Neural Information Processing*, pages 108–116. Springer.
- Scott, D. W. (1992). Multivariate density estimation: theory, practice, and visualization. John Wiley & Sons, Hoboken, NJ.
- Scott, D. W. (2015). Multivariate density estimation: theory, practice, and visualization. John Wiley & Sons, Hoboken, NJ.
- Sheather, S. J. (2004). Density estimation. Statistical Science, 19(4):588–597.
- Silverman, B. W. (1986). Density Estimation for Statistics and Data Analysis. Chapman and Hall, London, United Kingdom.

- Tarter, M. E. and Lock, M. D. (1993). *Model-free curve estimation*, volume 56. CRC Press, Boca Raton, FL.
- Tibshirani, R., Walther, G., and Hastie, T. (2001). Estimating the number of clusters in a data set via the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 63(2):411–423.
- Wang, X., Chen, H., Shen, D., and Huang, H. (2017). Cognitive impairment prediction in alzheimer's disease with regularized modal regression. In *Advances in Neural Information Processing Systems*, pages 1447–1457.
- Watson, G. S. (1964). Smooth regression analysis. Sankhyā: The Indian Journal of Statistics, Series A, pages 359–372.
- Wu, C. J. (1983). On the convergence properties of the em algorithm. *The Annals of Statistics*, 11(1):95–103.
- Yao, W. and Li, L. (2014). A new regression model: modal linear regression. *Scandinavian Journal of Statistics*, 41(3):656–671.
- Yao, W., Lindsay, B. G., and Li, R. (2012). Local modal regression. *Journal of nonparametric statistics*, 24(3):647–663.
- Zhou, H. and Huang, X. (2016). Nonparametric modal regression in the presence of measurement error. *Electronic Journal of Statistics*, 10(2):3579–3620.
- Zhou, H. and Huang, X. (2017). Bandwidth selection for nonparametric modal regression.

 To appear in Communications in Statistics Simulation and Computation.