

A Statistical Learning Approach to Modal Regression

Yunlong Feng

YLFENG@ALBANY.EDU

*Department of Mathematics and Statistics
State University of New York
The University at Albany
Albany, New York 12222, USA*

Jun Fan

JUNFAN@HKBU.EDU.HK

*Department of Mathematics
Hong Kong Baptist University
Kowloon, Hong Kong, China*

Johan A.K. Suykens

JOHAN.SUYKENS@ESAT.KULEUVEN.BE

*Department of Electrical Engineering
ESAT-STADIUS, KU Leuven
Kasteelpark Arenberg 10, Leuven
B-3001, Belgium*

Editor: Andreas Christmann

Abstract

This paper studies the nonparametric modal regression problem systematically from a statistical learning viewpoint. Originally motivated by pursuing a theoretical understanding of the maximum correntropy criterion based regression (MCCR), our study reveals that MCCR with a tending-to-zero scale parameter is essentially modal regression. We show that the nonparametric modal regression problem can be approached via the classical empirical risk minimization. Some efforts are then made to develop a framework for analyzing and implementing modal regression. For instance, the modal regression function is described, the modal regression risk is defined explicitly and its *Bayes* rule is characterized; for the sake of computational tractability, the surrogate modal regression risk, which is termed as the generalization risk in our study, is introduced. On the theoretical side, the excess modal regression risk, the excess generalization risk, the function estimation error, and the relations among the above three quantities are studied rigorously. It turns out that under mild conditions, function estimation consistency and convergence may be pursued in modal regression as in vanilla regression protocols such as mean regression, median regression, and quantile regression. On the practical side, the implementation issues of modal regression including the computational algorithm and the selection of the tuning parameters are discussed. Numerical validations on modal regression are also conducted to verify our findings.

Keywords: Nonparametric modal regression, empirical risk minimization, generalization bounds, kernel density estimation, statistical learning theory

1. Introduction

In this paper, we are interested in the nonparametric regression problem which aims at inferring the functional relation between input and output. Regression problems are concerned with the conditional distribution, which in practice can never be known in advance. Instead, normally, what one can access is only a set of observations drawn from the joint probability distribution. To state this problem mathematically, let us denote X as the explanatory variable that takes values in a compact metric space $\mathcal{X} \subset \mathbb{R}^d$ and Y that takes values in $\mathcal{Y} = \mathbb{R}$ as the response variable. Typically, we consider the following data-generating model

$$Y = f^*(X) + \epsilon,$$

where ϵ is the noise variable. In nonparametric regression problems, the purpose is to infer the unknown function f^* nonparametrically while certain assumptions on the noise variable ϵ may be imposed. As a compromise, regression estimators usually settle for learning a characterization of the conditional distribution by sifting information through observations generated above. Characterizations of the conditional distribution are versatile, where the several usual ones include the conditional mean, the conditional median, the conditional quantile, and the conditional mode. The versatility of the characterizations of the conditional distribution raises the question that which characterization we should pursue in regression problems. To answer this question, tremendous attention has been drawn in the statistics and machine learning communities. As a matter of fact, a significant part of parametric and nonparametric regression theory has been fostered to illuminate this question.

It is generally considered that each of the above-mentioned regression protocols has its own merits in its own regimes. For instance, it has been well understood that regression towards the conditional mean can be most effective if the noise is Gaussian or sub-Gaussian. Regression towards the conditional median or conditional quantile can be more robust in the absence of light-tailed noise or symmetric conditional distributions. In practice, the choice of the most appropriate regression protocol is usually decided by the type of data encountered. In the statistics and machine learning literature, these regression protocols have been studied extensively and understood well. In this study, we focus on a regression problem that has not been well studied in the statistical learning literature, namely, modal regression.

1.1. Modal Regression

Modal regression approaches the unknown truth f^* by regressing towards the *conditional mode function*. For a set of observations, the mode is the value that appears most frequently. While for a continuous random variable, the mode is the value at which its density function attains its peak value. The conditional mode function is denoted pointwisely as the mode of the conditional density of the dependent variable conditioned on the independent variable.

Previously proposed in Sager and Thisted (1982); Collomb et al. (1987) and studied in, e.g., Lee (1989, 1993), it is shown that one of the most appealing features of modal regression lies in its robustness to outliers, heavy-tailed noise, and skewed noise. Moreover, regression towards the conditional mode in some cases can be a better option when predicting the trends of observations. This is also the case in some real-world applications, as illustrated

in Matzner-Løfber et al. (1998); Einbeck and Tutz (2006); Yu et al. (2014). However, it seems to us that so far not enough attention has been given to the theory and applications of modal regression, especially in the statistical learning literature. As yet another regression protocol, the above-mentioned merits of modal regression suggest that it deserves far more attention than it has received, especially in the big data era today. This motivates our study on modal regression in this paper.

1.2. Historical Notes on Modal Regression

Modal regression is concerned with the mode. Studies on the mode estimation date back to the 1960s since the seminal work of Parzen (1962). It opens the door for kernel density estimation by proposing the *Parzen window* method, with the help of which the estimation of the mode can typically proceed. Many subsequent studies concerning theoretical as well as practical estimation of the mode have been emerging since then, among them are Chernoff (1964); Robertson and Cryer (1974); Fukunaga and Hostetler (1975); Eddy (1980); Comaniciu and Meer (2002), and Dasgupta and Kpotufe (2014).

In a regression setup, the concern of the conditional mode estimate gives birth to modal regression. As far as we are aware, the idea of regression towards the conditional mode was first proposed in Sager and Thisted (1982) in an isotonic regression setup. It was then specifically investigated in Collomb et al. (1987) when dealing with dependent observations. As a theoretical study, the main conclusion drawn there was the uniform convergence of the nonparametric mode estimator to the conditional mode function. Lately, in Lee (1989, 1993), some pioneering studies of modal regression were conducted. The tractability problem of mode regression was first discussed in their studies from, say, a supervised learning and risk minimization viewpoint. By considering some specific modal regression kernels, and assuming the existence of a global conditional mode function under a linear model assumption, they established the asymptotic normality of the resulting estimator. More and more attention to the theory and applications of modal regression has been attracted since the work in Yao et al. (2012); Yao and Li (2014) and Kemp and Santos Silva (2012). In Yao and Li (2014), a global mode was assumed to exist and take a linear form. Under proper assumptions on the conditional density of the noise variable, the implementation issues and the asymptotic normality of the estimator, as well as its robustness were explored. Recently, Chen et al. (2016b) presented an interesting study towards modal regression in which the conditional mode was sought by estimating the maximum of a joint density. By assuming a factorizable modal manifold collection, results on asymptotic error bounds as well as techniques for constructing confidence sets and prediction sets were provided.

To further disentangle the literature on modal regression, we can roughly categorize existing studies by tracing the thread of global or local approaches that they follow. For local approaches, the conditional mode is sought via maximizing a conditional density or a joint density which is typically estimated non-parametrically, e.g., by using kernel density estimators. Studies in Collomb et al. (1987); Samanta and Thavaneswaran (1990); Quintela-Del-Rio and Vieu (1997); Ould-Saïd (1997); Herrmann and Ziegler (2004); Ferraty et al. (2005); Gannoun et al. (2010); Yao et al. (2012); Chen et al. (2016b); Sasaki et al. (2016); Zhou and Huang (2016); Yao and Xiang (2016); Zhou and Huang (2019) fall into this category. For global approaches, the conditional mode is usually sought by maximizing

the kernel density estimator for the variable induced by the residual and assuming that the global mode is unique and belongs to a certain hypothesis space. To name a few, studies in Lee (1989, 1993); Lee and Kim (1998); Yao and Li (2014); Kemp and Santos Silva (2012); Baldauf and Santos Silva (2012); Yu and Aristodemou (2012); Lv et al. (2014); Salah and Françoise (2016) follow this line. It should be noticed that most studies based upon global approaches assume the existence (and also the uniqueness) of a global conditional mode function that is of a parametric form. While for the studies based upon local approaches, usually only the uniqueness assumption of the conditional mode function is imposed. Loosely speaking, modal regression estimators of the former case are nonparametric, while (semi-) parametric in the latter case.

Most of the above-mentioned studies are theoretical in nature. It should be noted that some application-oriented studies on modal regression have also been conducted. Among them, Matzner-Løfber et al. (1998) carried out an empirical comparison among three regression schemes, namely, the conditional mean regression, the conditional median regression, and the conditional mode regression, in nonparametric forecasting problems. They empirically observed that for certain datasets, e.g., the Old Faithful eruption prediction dataset, the mode can be a better option in forecasting than the mean and the median; Yu et al. (2014) discussed the mode-based regression problem in the big data context. Based on empirical evaluations on the Health Survey for England dataset, they argued that the mode could be an effective alternative for pattern-finding; Einbeck and Tutz (2006) dealt with the speed-flow data in traffic engineering by applying a multi-modal regression model.

1.3. Objectives of This Study and Our Contributions

As mentioned above, in the statistics literature, there exist some interesting studies towards modal regression from both theoretical and practical viewpoints. However, we notice that several problems related to the theoretical understanding as well as the practical implementations of modal regression remain unclear. For example:

- Modal regression regresses towards the conditional mode function, a direct estimation of which involves the estimation of a conditional or joint density. In fact, many of the existing studies on modal regression follow this approach. Notice that the explanatory variable may be high-dimensional vector-valued, which may make the estimation of the conditional or the joint density infeasible. This poses an important question: how to carry out modal regression without involving the estimation of a density function in a (possibly) high-dimensional space? According to the existing studies on modal regression, assuming the existence of a global conditional mode function and imposing some prior structure assumptions on it seem to be promising in avoiding estimating such a density. However, most existing studies of this type assume that the conditional mode function possesses a certain linear or parametric form. This could be restrictive in certain circumstances.
- With a modal regression estimator at hand, how can we evaluate its statistical performance? That is, how can we measure the approximation ability of the modal regression estimator to the conditional mode function? This concern is of great importance in nonparametric statistics as well as in machine learning as it is closely

related to the prediction ability of the estimator on future observations. On the other hand, concerning the implementation issues of modal regression, how can we perform model selection in modal regression?

To address the above two problems raised in modal regression, in this study, we propose to perform modal regression through the classical empirical risk minimization (ERM) scheme. Within the statistical learning framework, we then develop a learning theory framework for assessing the performance of the resulting modal regression estimator. Our contributions made in this study can be summarized as follows:

- The first main contribution of our study is that we present the first systematic statistical learning treatment on modal regression. This purpose is achieved by developing a statistical learning setup for modal regression, adapting it into the classical ERM framework, and conducting a learning theory analysis for modal regression estimators. The statistical learning approach to modal regression in this paper distinguishes our work from previous studies.
- The second main contribution of this study lies in that we develop a statistical learning framework for modal regression. To this end, the modal regression risk is devised, the *Bayes* rule of the modal regression risk is characterized, computationally tractable surrogates of the modal regression risk are introduced, and ERM schemes for modal regression are formulated.
- Following the ERM scheme, by assuming the existence of a global conditional mode function, the modal regression estimator in our study is pursued by maximizing a one-dimensional density estimator. This is more computationally tractable compared with the approaches adopted in most of the existing studies, in which the estimation of a possibly high-dimensional density is involved, as detailed in Section 3.5. This gives the third main contribution of this study.
- Another contribution made in this paper is that we present a learning theory analysis on the modal regression estimator resulted from the ERM scheme. The theoretical results in our analysis are concerned with the modal regression risk consistency, the generalization risk consistency, the function estimation ability of the modal regression estimator, and their relations, see Section 3 for details.
- It should be highlighted that, as we shall also explain below, the study in this paper is originally motivated by pursuing some further understanding of the maximum correntropy based regression (MCCR), which was recently investigated in Feng et al. (2015). In particular, this study is started with the realization that MCCR with a tending-to-zero scale parameter is modal regression, see Section 4 for details. It turns out that the study conducted in this paper brings us some new perspectives and a deeper understanding of MCCR.

1.4. Structure of This Paper

This paper is organized as follows: in Section 2, we formulate the modal regression problem within the statistical learning framework. To this end, we introduce the modal regres-

notation	meaning
\mathcal{X}, \mathcal{Y}	the independent variable space and the dependent variable space, respectively
X, Y	random variables taking values in \mathcal{X} and \mathcal{Y} , respectively
x, y	realizations of X and Y , respectively
\mathcal{M}	the function set comprised of all measurable function from \mathcal{X} to \mathbb{R}
ϵ	the noise variable specified by the residual $Y - f^*(X)$
\mathbf{z}	a set of n -size realizations of (X, Y) with $\mathbf{z} := \{(x_i, y_i)\}_{i=1}^n$
E_f	the random variable induced by the residual $Y - f(X)$
\mathcal{H}	a hypothesis space that is assumed to be a compact subset of $C(\mathcal{X})$
K_σ	a smoothing kernel with the bandwidth σ
ρ	the joint probability distribution of $X \times Y$
ρ_X	the marginal distribution of X
$L^2_{\rho_X}$	the function space of square-integrable functions with respect to ρ_X
p_{E_f} or p_f	the density function of the random variable E_f
$p_{Y X}$	the conditional density of Y conditioned on X
$p_{X,Y}$	the joint density of X and Y
$p_{\epsilon X}$	the conditional density of ϵ conditioned on X
f^*	the underlying truth function in modal regression, see formula (2.1)
f_M	the modal regression function or the conditional mode function, see formula (2.2)
$f_{\mathbf{z},\sigma}$	the empirical modal regression estimator in \mathcal{H} , see formula (2.4)
$f_{\mathcal{H},\sigma}$	the data-free modal regression estimator in \mathcal{H} , see formula (2.5)
$f_{\mathcal{H}}$	the data-free least squares regression estimator in \mathcal{H}
$\mathcal{R}(f)$	the modal regression risk for the hypothesis $f : \mathcal{X} \rightarrow \mathbb{R}$
$\mathcal{R}^\sigma(f)$	the data-free generalization risk for the hypothesis $f : \mathcal{X} \rightarrow \mathbb{R}$
$\mathcal{R}_n^\sigma(f)$	the empirical generalization risk for the hypothesis $f : \mathcal{X} \rightarrow \mathbb{R}$

Table 1: A list of notations and their definitions in this paper

sion function in Subsection 2.1. We define the modal regression risk and characterize its *Bayes* rule in Subsection 2.2. A kernel density estimation interpretation and an empirical risk minimization perspective of modal regression are provided in Subsections 2.3 and 2.4, respectively. Section 3 is devoted to developing a learning theory for modal regression. The modal regression calibration problem (see Subsection 3.2), the convergence of the excess generalization risk (see Subsection 3.3), and the function estimation calibration problem (see Subsection 3.4) are studied by applying standard learning theory arguments. Comparisons between our study and the existing ones are also mentioned in this section. In Section 4, we interpret MCCR from a modal regression viewpoint by suggesting that MCCR with a tending-to-zero scale parameter is essentially modal regression. Since one of the main motivations of the present study is to understand MCCR within the statistical learning framework and having realized that MCCR with a tending-to-zero scale parameter is modal regression, we, therefore, retrospect MCCR in Section 4.2 by applying the theory developed in Section 3 and depict a general picture of MCCR. Section 5 is concerned with the implementation issues in modal regression such as model selection and computational

algorithms. Numerical validations will be provided in this section. We close this paper in Section 6 with conclusions. For the sake of readability, a list of notations and their definitions in this paper is provided in Table 1.

2. A Statistical Learning Framework for Modal Regression

2.1. Formulating the Modal Regression Problem

We first formulate the modal regression problem formally in this subsection. To this end, we first assume that we are given a set of i.i.d observations \mathbf{z} that are generated by

$$Y = f^*(X) + \epsilon, \tag{2.1}$$

where the mode of the conditional distribution of ϵ at any $x \in \mathcal{X}$ is assumed to be zero. That is, $\text{mode}(\epsilon | X = x) := \arg \max_{t \in \mathbb{R}} p_{\epsilon|X}(t | X = x) = 0$ for any $x \in \mathcal{X}$, where $p_{\epsilon|X}$ is the conditional density of ϵ conditioned on X . It is obvious from (2.1) that under the zero-mode noise assumption, it holds that $\text{mode}(Y|X) = f^*(X)$. We further assume that $p_{\epsilon|X}$ is continuous and bounded on \mathbb{R} for any $x \in \mathcal{X}$. Here, it should be remarked that in this study we do not assume either the homogeneity or the symmetry of the distribution of the noise ϵ . In other words, the heterogeneity of the distribution of the residuals or the skewed noise distribution is allowed.

In modal regression problems, we aim at approximating the *modal regression function* (see formula (1.1) in Collomb et al. (1987)):

Definition 1 (Modal Regression Function) *The modal regression function $f_M : \mathcal{X} \rightarrow \mathbb{R}$ is defined as*

$$f_M(x) := \arg \max_{t \in \mathbb{R}} p_{Y|X}(t | X = x), \quad x \in \mathcal{X}, \tag{2.2}$$

where $p_{Y|X}(\cdot|X)$ denotes the conditional density of Y conditioned on X .

Throughout this paper, we assume that the modal regression function f_M is well-defined on \mathcal{X} . That is, $\arg \max_{t \in \mathbb{R}} p_{Y|X}(t | X = x)$ is assumed to exist and be unique for any fixed $x \in \mathcal{X}$. Obviously, this is equivalent to assuming the existence and uniqueness of the global mode of the conditional density $p_{Y|X}$. On the other hand, due to the zero-mode assumption of the conditional distribution of ϵ in (2.1) for any $x \in \mathcal{X}$, we know that $f_M \equiv f^*$. Consequently, the learning for modal regression problem is equivalent to the problem of learning the modal regression function f_M , and thus f^* . Said differently, f_M is the so-called target hypothesis.

From the definition, the modal regression function f_M is defined as the maximum of the conditional density $p_{Y|X}$ conditioned on X . Note that, maximizing the conditional density is equivalent to maximizing the joint density $p_{X,Y}$ for any fixed realization of X . Therefore, it is direct to see that one can approximate f_M by maximizing the conditional density $p_{Y|X}$ or the joint density $p_{X,Y}$, both of which can be estimated via kernel density estimation. This is, in fact, what most of the existing studies on modal regression do (see e.g., Collomb et al., 1987; Chen et al., 2016b; Yao and Xiang, 2016). However, estimating the conditional density $p_{Y|X}$ or the joint density $p_{X,Y}$ via kernel density estimation suffers

from the *curse of dimensionality* and is not feasible when the dimension of the input space is high. In this study, we are interested in an empirical risk minimization approach that is dimension-insensitive as formulated later.

2.2. Modeling the Modal Regression Risk and Characterizing the *Bayes* Rule

To be in a position to carry out a statistical learning assessment of modal regression, besides the target hypothesis defined above, we also need to devise a fitting risk that measures the goodness-of-fit when a candidate hypothesis is considered. The newly devised fitting risk should vote the target hypothesis (2.2) as the best candidate when the hypothesis space is sufficiently large. This gives the main purpose of this subsection.

Definition 2 (Modal Regression Risk) For a measurable function $f : \mathcal{X} \rightarrow \mathbb{R}$, its *modal regression risk* $\mathcal{R}(f)$ is defined as

$$\mathcal{R}(f) = \int_{\mathcal{X}} p_{Y|X}(f(x)|X=x) d\rho_{\mathcal{X}}(x). \quad (2.3)$$

Analogously to learning for regression and classification scenarios (see, e.g., Cucker and Zhou, 2007; Steinwart and Christmann, 2008), we denote the *Bayes rule* of modal regression as the “best” hypothesis favored by the above modal regression risk over the measurable function set \mathcal{M} (comprised of all measurable functions from \mathcal{X} to \mathbb{R}). The following conclusion indicates that the target hypothesis f_M is exactly the *Bayes* rule of modal regression.

Theorem 3 The modal regression function f_M in (2.2) gives the *Bayes* rule of modal regression. That is,

$$f_M = \arg \max_{f \in \mathcal{M}} \mathcal{R}(f).$$

Proof Recall that the conditional mode function f_M is given as

$$f_M(x) = \arg \max_{t \in \mathbb{R}} p_{Y|X}(t|X=x), \quad x \in \mathcal{X}.$$

Following the modal regression risk defined in Definition 2, for any measurable function $f \in \mathcal{M}$, we have

$$\mathcal{R}(f) = \int_{\mathcal{X}} p_{Y|X}(f(x)|X=x) d\rho_{\mathcal{X}}(x) \leq \int_{\mathcal{X}} p_{Y|X}(f_M(x)|X=x) d\rho_{\mathcal{X}}(x) = \mathcal{R}(f_M),$$

which directly yields

$$f_M = \arg \max_{f \in \mathcal{M}} \mathcal{R}(f).$$

This completes the proof of Theorem 3. ■

The plausibility of the above-defined modal regression risk stems from the fact that f_M is the *Bayes* rule of modal regression, as justified by Theorem 3. With the modal regression

risk being defined and recalling that f_M maximizes the modal regression risk, the most direct way to learn f_M is to maximize the sample analogy of the modal regression risk. Unfortunately, this is intractable since the discretization of an unknown conditional density is involved. In the next subsection, to circumvent this problem, we introduce a surrogate of the modal regression risk.

Remark 4 *We now give a remark on the terminology “risk”. For any measurable function $f : \mathcal{X} \rightarrow \mathbb{R}$, the modal regression risk $\mathcal{R}(f)$ in Definition 2 can be regarded as a measure of the extent to which the function f fits the Bayes rule f_M in the $\mathcal{R}(\cdot)$ sense. Therefore, the terminology “risk” is not used as what is commonly referred to in the statistical learning literature. However, in what follows, given the one-to-one correspondence between the corresponding maximization and minimization problems, we still term $\mathcal{R}(f)$ as the (modal regression) risk of f .*

2.3. Learning for Modal Regression via Kernel Density Estimation

We now show that the modal regression problem can be tackled by applying the kernel density estimation technique. To this purpose, let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a measurable function and denote E_f as the random variable induced by the residual $Y - f(X)$, where the subscript f indicates its dependence on f . We also denote p_{E_f} , or simply p_f , as the density function of the random variable E_f and denote $p_{\epsilon|X}$ as the conditional density of the random variable $\epsilon = Y - f^*(X)$. The following theorem, which was first established in Fan et al. (2016), relates the modal regression risk of f to $p_{\epsilon|X}$ and p_{E_f} .

Theorem 5 *Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a measurable function. Then,*

$$\int_{\mathcal{X}} p_{\epsilon|X}(\cdot + f(x) - f^*(x)|X = x) d\rho_{\mathcal{X}}(x)$$

is a density of the random variable $E_f := Y - f(X)$, which is denoted as p_{E_f} . Correspondingly, we have $p_{E_f}(0) = \mathcal{R}(f)$.

Proof From the model assumption that $\epsilon = Y - f^*(X)$, we have

$$\epsilon = E_f + f(X) - f^*(X).$$

As a result, the density function of the error variable E_f can be expressed as

$$\int_{\mathcal{X}} p_{\epsilon|X}(\cdot + f(x) - f^*(x)|X = x) d\rho_{\mathcal{X}}(x)$$

and denoted by p_{E_f} . Moreover, from the definition of the risk functional $\mathcal{R}(\cdot)$ in (2.3), we know that

$$\begin{aligned} p_{E_f}(0) &= \int_{\mathcal{X}} p_{\epsilon|X}(f(x) - f^*(x)|X = x) d\rho_{\mathcal{X}}(x) \\ &= \int_{\mathcal{X}} p_{Y|X}(f(x)|X = x) d\rho_{\mathcal{X}}(x) \\ &= \mathcal{R}(f). \end{aligned}$$

This completes the proof of Theorem 5. ■

From Theorem 5, the hypothesis f that maximizes the modal regression risk $\mathcal{R}(f)$ is the one that maximizes the density of $E_f := Y - f(X)$ at 0, which can be estimated non-parametrically. In this study, the kernel density estimation technique is tailored to modal regression with the help of the modal regression kernel defined below.

Definition 6 (Modal Regression Kernel) *A kernel $K_\sigma : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$ is said to be a **modal regression kernel** with the representing function ϕ and the bandwidth parameter $\sigma > 0$ if there exists a function $\phi : \mathbb{R} \rightarrow [0, \infty)$ such that $K_\sigma(u_1, u_2) = \phi\left(\frac{u_1 - u_2}{\sigma}\right)$ for any $u_1, u_2 \in \mathbb{R}$, $\phi(u) = \phi(-u)$, $\phi(u) \leq \phi(0)$ for any $u \in \mathbb{R}$, and $\int_{\mathbb{R}} \phi(u) du = 1$.*

According to Definition 6, it is easy to see that common smoothing kernels (see, e.g., Wand and Jones, 1994) such as the Naive kernel, the Gaussian kernel, the Epanechnikov kernel, and the Triangular kernel are modal regression kernels. Their corresponding representing functions can be easily deduced with simple computations. For a modal regression kernel K_σ with the representing function ϕ , throughout this paper, without loss of generality, we assume $\phi(0) = 1$.

As a consequence of Theorem 5, for any measurable function f , we know that $p_f(0) = \mathcal{R}(f)$. With the help of a modal regression kernel K_σ , it is immediate to see that an empirical kernel density estimator \hat{p}_f for p_f at 0 can be formulated as follows

$$\hat{p}_f(0) = \frac{1}{n\sigma} \sum_{i=1}^n K_\sigma(y_i - f(x_i), 0) = \frac{1}{n\sigma} \sum_{i=1}^n K_\sigma(y_i, f(x_i)) := \mathcal{R}_n^\sigma(f).$$

Therefore, when confined to a hypothesis space \mathcal{H} , learning a function f that maximizes the modal regression risk is cast as learning the function f that maximizes the value of the empirical density estimator \hat{p}_f at 0. Thus, the empirical target hypothesis is modeled as

$$\begin{aligned} f_{\mathbf{z}, \sigma} &:= \arg \max_{f \in \mathcal{H}} \hat{p}_f(0) \\ &= \arg \max_{f \in \mathcal{H}} \mathcal{R}_n^\sigma(f), \end{aligned} \tag{2.4}$$

where \mathcal{H} is assumed to be a compact subset of $C(\mathcal{X})$ throughout this paper. The population version of $f_{\mathbf{z}, \sigma}$ can be expressed as

$$f_{\mathcal{H}, \sigma} := \arg \max_{f \in \mathcal{H}} \mathcal{R}^\sigma(f), \tag{2.5}$$

where $\mathcal{R}^\sigma(\cdot)$ is the expectation of $\mathcal{R}_n^\sigma(\cdot)$ with respect to the random samples \mathbf{z} and for any $f : \mathcal{X} \rightarrow \mathbb{R}$, it can be expressed as

$$\mathcal{R}^\sigma(f) = \frac{1}{\sigma} \int_{\mathcal{X} \times \mathcal{Y}} \phi\left(\frac{y - f(x)}{\sigma}\right) d\rho(x, y).$$

The risk functional $\mathcal{R}^\sigma(f)$ defined above gives the **generalization risk** of f when a modal regression kernel K_σ with the representing function ϕ is adopted. As we shall see later, it can be seen as a surrogate of the true modal regression risk $\mathcal{R}(f)$ since $\mathcal{R}^\sigma(f)$ approximates $\mathcal{R}(f)$ when $\sigma \rightarrow 0$. The interpretation of modal regression from a kernel density estimation viewpoint explains the requirement that $\int_{\mathbb{R}} \phi(u) du = 1$ in Definition 6.

2.4. Modal Regression: an Empirical Risk Minimization View

In the preceding subsection, we showed that the modal regression scheme (2.4) can be interpreted from a kernel density estimation point of view. Maximizing the value of the kernel density estimator for E_f at 0 encourages the considered hypothesis f to approximate the projection of the *Bayes* rule onto \mathcal{H} , i.e., $f_{\mathcal{H},\sigma}$. In this subsection, we show that one can also interpret the modal regression scheme (2.4) by using the language of empirical risk minimization.

To proceed, let us consider a modal regression kernel K_σ with the representing function ϕ and the scale parameter $\sigma > 0$. We then introduce the following distance-based modal regression loss $\phi_\sigma : \mathbb{R} \rightarrow [0, \infty)$:

$$\phi_\sigma(y - f(x)) = \sigma^{-1} (1 - \phi((y - f(x))\sigma^{-1})). \quad (2.6)$$

Based on the newly introduced loss ϕ_σ , the modal regression scheme (2.4) can be reformulated as follows

$$f_{\mathbf{z},\sigma} = \arg \min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \phi_\sigma(y_i, f(x_i)), \quad (2.7)$$

and, similarly, its data-free counterpart can be formulated as

$$f_{\mathcal{H},\sigma} = \arg \min_{f \in \mathcal{H}} \int_{\mathcal{X} \times \mathcal{Y}} \phi_\sigma(y, f(x)) d\rho. \quad (2.8)$$

It is easy to see that the empirical estimator (2.7) is an M-estimator and the two formulations of $f_{\mathbf{z},\sigma}$ in (2.4) and (2.7) are, in fact, equivalent. Similarly, one also obtains the same target hypothesis from (2.5) and (2.8).

Remark 7 *For formulation simplification, whenever referred to herein, $f_{\mathbf{z},\sigma}$ and $f_{\mathcal{H},\sigma}$ will be pointed to the estimators formulated by (2.4) and (2.5), respectively, while keeping in mind that the conducted analysis on $f_{\mathbf{z},\sigma}$ is inspired by and within the ERM framework.*

3. A Learning Theory of Modal Regression

In this section, we aim to develop a learning theory for modal regression which can be used to assess the statistical learning performance of the modal regression estimator $f_{\mathbf{z},\sigma}$.

3.1. Learning the Conditional Mode: Three Building Blocks

In Section 2, for a given hypothesis f , the modal regression risk $\mathcal{R}(f)$ is defined; moreover, it turns out that f_M is the *Bayes* rule of modal regression. On the other hand, we show that the modal regression estimator can be learned via maximizing the risk functional $\mathcal{R}_n^\sigma(\cdot)$. Recalling that the central concern in learning theory is risk consistency under various notions and following the clue of existing learning theory studies on the binary-classification problem, it is natural and necessary to investigate the following three problems:

1. The problem of the excess generalization risk consistency and convergence rates, i.e., the convergence from $\mathcal{R}^\sigma(f_{\mathbf{z},\sigma})$ to $\mathcal{R}^\sigma(f^*)$.

2. The modal regression calibration problem, i.e., whether the convergence from $\mathcal{R}^\sigma(f_{\mathbf{z},\sigma})$ to $\mathcal{R}^\sigma(f^*)$ implies the convergence from $\mathcal{R}(f_{\mathbf{z},\sigma})$ to $\mathcal{R}(f^*)$?
3. The function estimation calibration problem, i.e., whether the convergence from $\mathcal{R}(f_{\mathbf{z},\sigma})$ to $\mathcal{R}(f^*)$ implies the convergence from $f_{\mathbf{z},\sigma}$ to f^* ?

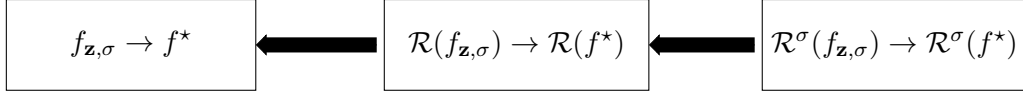


Figure 1: An illustration of the three building blocks in learning for modal regression. The left block stands for the function estimation consistency of $f_{\mathbf{z},\sigma}$, the middle block denotes the modal regression consistency of $f_{\mathbf{z},\sigma}$, while the right block represents the excess generalization risk consistency of $f_{\mathbf{z},\sigma}$.

The above three problems are fundamental in conducting a learning theory analysis on modal regression and serve as three main building blocks. Detailed explorations will be expanded in the following subsections.

3.2. Towards the Modal Regression Calibration Problem

We first investigate the modal regression calibration problem stated in Question 1, i.e., whether the convergence from $\mathcal{R}^\sigma(f_{\mathbf{z},\sigma})$ to $\mathcal{R}^\sigma(f^*)$ implies the convergence from $\mathcal{R}(f_{\mathbf{z},\sigma})$ to $\mathcal{R}(f^*)$. To this end, we need to confine ourselves to the calibrated modal regression kernel defined below.

Definition 8 (Calibrated Modal Regression Kernel) *A modal regression kernel K_σ with the representing function ϕ is said to be a **calibrated modal regression kernel** if it satisfies the following conditions:*

- (i) ϕ is bounded;
- (ii) ϕ is Lipschitz continuous on \mathbb{R} with the Lipschitz constant L ;
- (iii) $\int_{\mathbb{R}} u^2 \phi(u) du < \infty$.

Another restriction we need to impose is on the conditional density $p_{\epsilon|X}$ as follows:

Assumption 1 *The conditional density of ϵ given X , namely, $p_{\epsilon|X}$, is second-order continuously differentiable and $\|p''_{\epsilon|X}\|_\infty$ is bounded from above.*

Theorem 9 *Suppose that Assumption 1 holds and let K_σ be a calibrated modal regression kernel with the representing function ϕ and the scale parameter σ . For any measurable function $f : \mathcal{X} \rightarrow \mathbb{R}$, it holds that*

$$\left| \{\mathcal{R}(f^*) - \mathcal{R}(f)\} - \{\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f)\} \right| \leq c_1 \sigma^2,$$

where $c_1 = \|p''_{\epsilon|X}\|_\infty \int_{\mathbb{R}} u^2 \phi(u) du$.

Proof Recalling the definition of the risk functional $\mathcal{R}^\sigma(f)$ for any measurable function $f : \mathcal{X} \rightarrow \mathbb{R}$ and applying Taylor's Theorem to the conditional density $p_{\epsilon|X}$, we have

$$\begin{aligned}
 \mathcal{R}^\sigma(f) &= \frac{1}{\sigma} \int_{\mathcal{X} \times \mathcal{Y}} \phi\left(\frac{y - f(x)}{\sigma}\right) d\rho(x, y) \\
 &= \frac{1}{\sigma} \int_{\mathcal{X}} \int_{\mathbb{R}} \phi\left(\frac{t - (f(x) - f^*(x))}{\sigma}\right) p_{\epsilon|X}(t | X = x) dt d\rho_{\mathcal{X}}(x) \\
 &= \int_{\mathcal{X}} \int_{\mathbb{R}} \phi(u) p_{\epsilon|X}(f(x) - f^*(x) + \sigma u | X = x) du d\rho_{\mathcal{X}}(x) \\
 &= \int_{\mathcal{X}} \int_{\mathbb{R}} \phi(u) p_{\epsilon|X}(f(x) - f^*(x) | X = x) du d\rho_{\mathcal{X}}(x) \\
 &\quad + \sigma \int_{\mathcal{X}} \int_{\mathbb{R}} u \phi(u) p'_{\epsilon|X}(f(x) - f^*(x) | X = x) du d\rho_{\mathcal{X}}(x) \\
 &\quad + \frac{\sigma^2}{2} \int_{\mathcal{X}} \int_{\mathbb{R}} u^2 \phi(u) p''_{\epsilon|X}(\eta_x | X = x) du d\rho_{\mathcal{X}}(x),
 \end{aligned} \tag{3.1}$$

where, for any fixed $x \in \mathcal{X}$, the point η_x lies between $f(x) - f^*(x)$ and $f(x) - f^*(x) + \sigma u$.

The fact that K_σ is a calibrated modal regression kernel with the representing function ϕ ensures $\int_{\mathbb{R}} \phi(u) du = 1$ and reminds the symmetry of ϕ on \mathbb{R} , which further indicates that $\int_{\mathbb{R}} u \phi(u) du = 0$. On the other hand, the fact that

$$\mathcal{R}(f) = \int_{\mathcal{X}} p_{\epsilon|X}(f(x) - f^*(x) | X = x) d\rho_{\mathcal{X}}(x),$$

together with Equalities (3.1) yields

$$|\mathcal{R}^\sigma(f) - \mathcal{R}(f)| \leq \frac{\sigma^2}{2} \left(\|p''_{\epsilon|X}\|_\infty \int_{\mathbb{R}} u^2 \phi(u) du \right).$$

Denoting $c_1 := \|p''_{\epsilon|X}\|_\infty \int_{\mathbb{R}} u^2 \phi(u) du$, we accomplish the proof of Theorem 9. \blacksquare

Remark 10 *The proof of Theorem 9 indicates that $\mathcal{R}^\sigma(f)$ is a second-order approximation (with respect to σ) of $\mathcal{R}(f)$ since $\mathcal{R}^\sigma(f) - \mathcal{R}(f) = \mathcal{O}(\sigma^2)$. In fact, if a higher-order kernel (see e.g., Section 2.8 in Wand and Jones, 1994) is used, a higher-order approximation of $\mathcal{R}(f)$ can be expected.*

From the proof of Theorem 9, we see that when K_σ is a calibrated modal regression kernel with the representing function ϕ and the scale parameter σ , for any measurable function $f : \mathcal{X} \rightarrow \mathbb{R}$, the generalization risk $\mathcal{R}^\sigma(f)$ approaches the true modal regression risk $\mathcal{R}(f)$ provided that $\sigma \rightarrow 0$. Therefore, in the above sense, $\mathcal{R}^\sigma(f)$ can be considered as a relaxation of $\mathcal{R}(f)$. On the other hand, Theorem 9 indicates that the difference between the **excess modal regression risk** $\mathcal{R}(f^*) - \mathcal{R}(f)$ and the **excess generalization risk** $\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f)$ can be upper bounded by $\mathcal{O}(\sigma^2)$. Clearly, under the assumptions of Theorem 9, when $\sigma \rightarrow 0$, $\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f)$ also approaches $\mathcal{R}(f^*) - \mathcal{R}(f)$. In this sense, Theorem 9 establishes a *comparison theorem* akin to the one in the classification scenario (see Zhang, 2004; Bartlett et al., 2006). This elucidates the terminology—the calibrated modal regression kernel, and the terminology—the modal regression calibration problem.

3.3. Towards the Convergence Rates of the Excess Generalization Risk

One of the main focuses in learning theory is the generalization ability of a learning algorithm that measures its out-of-sample prediction ability. It plays an important role in designing learning algorithms with theoretical guarantees. In this subsection, we derive the generalization bounds for the modal regression estimator $f_{\mathbf{z},\sigma}$, i.e., the convergence rates of $\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma})$, by means of learning theory arguments. The following assumption is needed for this purpose:

Assumption 2 *We make the following assumptions:*

- (i) *There exists a positive constant M such that $\|f^*\|_\infty \leq M$;*
- (ii) $\sup_{t \in \mathbb{R}, x \in \mathcal{X}} p_{\epsilon|X}(t | X = x) = c_2 < \infty$;
- (iii) *For any $\varepsilon > 0$, there exists an exponent p with $0 < p < 2$ such that the ℓ^2 -empirical covering number (with radius ε) of \mathcal{H} , denoted as $\mathcal{N}_{2,\mathbf{x}}(\mathcal{H}, \varepsilon)$, satisfies*

$$\log \mathcal{N}_{2,\mathbf{x}}(\mathcal{H}, \varepsilon) \lesssim \varepsilon^{-p},$$

where the definition of the empirical covering number is provided below (see also Anthony and Bartlett (2009)), and the notation $a \lesssim b$ for $a, b \in \mathbb{R}$ means that there exists a positive constant c such that $a \leq cb$.

Definition 11 (ℓ^2 -empirical Covering Number) *Let \mathcal{F} be a set of functions on \mathcal{X} and $\mathbf{x} = \{x_1, \dots, x_m\} \subset \mathcal{X}$. The metric $d_{2,\mathbf{x}}$ is defined on \mathcal{F} by*

$$d_{2,\mathbf{x}}(f, g) = \left\{ \frac{1}{m} \sum_{i=1}^m (f(x_i) - g(x_i))^2 \right\}^{1/2}.$$

For every $\varepsilon > 0$, the ℓ^2 -empirical covering number of \mathcal{F} with respect to $d_{2,\mathbf{x}}$ is defined as

$$\mathcal{N}_{2,\mathbf{x}}(\mathcal{F}, \varepsilon) = \inf \left\{ \ell \in \mathbb{N} : \exists \{f_i\}_{i=1}^\ell \text{ such that } \mathcal{F} = \cup_{i=1}^\ell \{f \in \mathcal{F} : d_{2,\mathbf{x}}(f, f_i) \leq \varepsilon\} \right\}.$$

Restrictions in Assumption 2 are fairly standard if we recall that the hypothesis space \mathcal{H} is assumed to be a compact subset of $C(\mathcal{X})$. In what follows, without loss of generality, we also assume that $\|f\|_\infty \leq M$ for any $f \in \mathcal{H}$. The following error decomposition lemma is helpful in bounding the excess generalization error.

Lemma 12 *Let $f_{\mathbf{z},\sigma}$ be produced by (2.4) and assume that $f^* \in \mathcal{H}$. Then we have*

$$\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma}) \leq \mathcal{R}^\sigma(f_{\mathcal{H},\sigma}) - \mathcal{R}_n^\sigma(f_{\mathcal{H},\sigma}) + \mathcal{R}_n^\sigma(f_{\mathbf{z},\sigma}) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma}).$$

Proof Recalling that $f_{\mathcal{H},\sigma} = \arg \max_{f \in \mathcal{H}} \mathcal{R}^\sigma(f)$, we have

$$\begin{aligned} \mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma}) &\leq \mathcal{R}^\sigma(f_{\mathcal{H},\sigma}) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma}) \\ &\leq \mathcal{R}^\sigma(f_{\mathcal{H},\sigma}) - \mathcal{R}_n^\sigma(f_{\mathcal{H},\sigma}) + \mathcal{R}_n^\sigma(f_{\mathcal{H},\sigma}) - \mathcal{R}_n^\sigma(f_{\mathbf{z},\sigma}) + \mathcal{R}_n^\sigma(f_{\mathbf{z},\sigma}) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma}) \\ &\leq \mathcal{R}^\sigma(f_{\mathcal{H},\sigma}) - \mathcal{R}_n^\sigma(f_{\mathcal{H},\sigma}) + \mathcal{R}_n^\sigma(f_{\mathbf{z},\sigma}) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma}), \end{aligned}$$

where the last inequality is due to the fact that the quantity $\mathcal{R}_n^\sigma(f_{\mathcal{H},\sigma}) - \mathcal{R}_n^\sigma(f_{\mathbf{z},\sigma})$ is at most zero. This completes the proof of Lemma 12. \blacksquare

The following lemma, established in Wu et al. (2007), provides a Bernstein-type concentration inequality for function-valued random variables. It was proved by applying the local Rademacher complexity arguments developed in Bartlett et al. (2005).

Lemma 13 *Let \mathcal{F} be a class of bounded measurable functions. Assume that there are constants $\gamma \in [0, 1]$ and $B, c_\gamma > 0$ such that $\|f\|_\infty \leq B$ and $\mathbb{E}f^2 \leq c_\gamma(\mathbb{E}f)^\gamma$ for every $f \in \mathcal{F}$. If for some $c_p > 0$ and $0 < p < 2$,*

$$\log \mathcal{N}_{2,\mathbf{x}}(\mathcal{F}, \varepsilon) \leq c_p \varepsilon^{-p}, \quad \forall \varepsilon > 0,$$

then there exists a constant c'_p depending only on p such that for any $t > 0$, with probability at least $1 - e^{-t}$, it holds that

$$\mathbb{E}f - \frac{1}{n} \sum_{i=1}^n f(z_i) \leq \frac{1}{2} \eta^{1-\gamma} (\mathbb{E}f)^\gamma + c'_p \eta + 2 \left(\frac{c_\gamma t}{n} \right)^{\frac{1}{2-\gamma}} + \frac{18Bt}{n}, \quad \forall f \in \mathcal{F},$$

where

$$\eta = \max \left\{ c_\gamma^{\frac{2-p}{4-2\gamma+p\gamma}} \left(\frac{c_p}{n} \right)^{\frac{2}{4-2\gamma+p\gamma}}, B^{\frac{2-p}{2+p}} \left(\frac{c_p}{n} \right)^{\frac{2}{2+p}} \right\}.$$

Theorem 14 *Suppose that Assumption 2 holds, $f^* \in \mathcal{H}$, and the risk functional $\mathcal{R}^\sigma(\cdot)$ is defined in association with a calibrated modal regression kernel K_σ and the representing function ϕ . Let $f_{\mathbf{z},\sigma}$ be produced by (2.4) with $\sigma \leq 1$. Then for any $0 < \delta < 1$, with probability at least $1 - \delta$, it holds that*

$$\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma}) \lesssim \left(\frac{1}{n\sigma} + \frac{\sigma^{-\frac{2+3p}{4}}}{n^{1/2}} + \frac{\sigma^{-\frac{2+3p}{2+p}}}{n^{\frac{2}{2+p}}} \right) \log \left(\frac{1}{\delta} \right).$$

Proof We prove the theorem by applying Lemma 13 to the following function-valued random variable on $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$:

$$\xi(z) := \frac{1}{\sigma} \phi \left(\frac{y - f_{\mathcal{H},\sigma}(x)}{\sigma} \right) - \frac{1}{\sigma} \phi \left(\frac{y - f(x)}{\sigma} \right), \quad (3.2)$$

where $f_{\mathcal{H},\sigma}$ is given in (2.5) and $f \in \mathcal{H}$. Due to the boundedness assumption of ϕ , it is easy to see that $|\xi(z)| \leq 2\|\phi\|_\infty/\sigma$. Moreover, recalling the definition of the risk functional $\mathcal{R}^\sigma(\cdot)$, the following inequality holds

$$\begin{aligned} \mathbb{E}\xi^2 &= \mathbb{E} \left[\frac{1}{\sigma} \phi \left(\frac{Y - f_{\mathcal{H},\sigma}(X)}{\sigma} \right) - \frac{1}{\sigma} \phi \left(\frac{Y - f(X)}{\sigma} \right) \right]^2 \\ &\leq \frac{2\|\phi\|_\infty}{\sigma} (\mathcal{R}^\sigma(f_{\mathcal{H},\sigma}) + \mathcal{R}^\sigma(f)). \end{aligned} \quad (3.3)$$

From the proof of Theorem 9, we know that

$$\mathcal{R}^\sigma(f_{\mathcal{H},\sigma}) \leq \mathcal{R}(f_{\mathcal{H},\sigma}) + \frac{\sigma^2}{2} \left(\|p''_{\epsilon|X}\|_\infty \int_{\mathbb{R}} u^2 \phi(u) du \right).$$

Similarly, we also have

$$\mathcal{R}^\sigma(f) \leq \mathcal{R}(f) + \frac{\sigma^2}{2} \left(\|p''_{\epsilon|X}\|_\infty \int_{\mathbb{R}} u^2 \phi(u) du \right).$$

The above two inequalities together with the bound for $\mathbb{E}\xi^2$ and the fact that $\sigma \leq 1$ yield

$$\begin{aligned} \mathbb{E}\xi^2 &\leq \frac{2\|\phi\|_\infty}{\sigma} (\mathcal{R}(f_{\mathcal{H},\sigma}) + \mathcal{R}(f) + c_1\sigma^2) \\ &\leq \frac{2\|\phi\|_\infty}{\sigma} (p_{f_{\mathcal{H},\sigma}}(0) + p_f(0) + c_1\sigma^2) \\ &\lesssim \sigma^{-1}, \end{aligned}$$

where the last inequality is due to the boundedness assumption of the conditional density of ϵ while the second inequality is a consequence of Theorem 5.

Recalling that ϕ is Lipschitz continuous on \mathbb{R} with the Lipschitz constant L , for any $f_1, f_2 \in \mathcal{H}$, we thus have

$$\left| \frac{1}{\sigma} \phi \left(\frac{y - f_1(x)}{\sigma} \right) - \frac{1}{\sigma} \phi \left(\frac{y - f_2(x)}{\sigma} \right) \right| \leq \frac{L}{\sigma^2} \|f_1 - f_2\|_\infty.$$

Consequently, if we denote $\mathcal{F}_{\mathcal{H}}$ as the following set

$$\mathcal{F}_{\mathcal{H}} := \left\{ g \mid g(z) = \frac{1}{\sigma} \phi \left(\frac{y - f_{\mathcal{H},\sigma}(x)}{\sigma} \right) - \frac{1}{\sigma} \phi \left(\frac{y - f(x)}{\sigma} \right), f \in \mathcal{H} \right\},$$

then Assumption 2 (iii) implies that

$$\log \mathcal{N}_{2,\mathbf{x}}(\mathcal{F}_{\mathcal{H}}, \varepsilon) \leq \log \mathcal{N}_{2,\mathbf{x}}(\mathcal{H}, \varepsilon\sigma^2/L) \lesssim (\varepsilon\sigma^2)^{-p}.$$

Applying Lemma 13 to the random variable ξ with $B = 2\|\phi\|_\infty/\sigma$, $\gamma = 0$, $c_p = \sigma^{-2p}$, and $c_\gamma = \sigma^{-1}$, then for any $0 < \delta < 1$, with probability at least $1 - \delta$, it holds that

$$\mathcal{R}^\sigma(f_{\mathcal{H},\sigma}) - \mathcal{R}^\sigma(f) - (\mathcal{R}_n^\sigma(f_{\mathcal{H},\sigma}) - \mathcal{R}_n^\sigma(f)) \lesssim \left(\frac{1}{n\sigma} + \frac{\sigma^{-\frac{2+3p}{4}}}{n^{1/2}} + \frac{\sigma^{-\frac{2+3p}{2+p}}}{n^{\frac{2}{2+p}}} \right) \log \left(\frac{1}{\delta} \right).$$

Noticing that the above inequality holds for any $f \in \mathcal{H}$ and recalling Lemma 12, we obtain the desired conclusion in Theorem 14. \blacksquare

The generalization bounds in Theorem 14 are derived for the case when the parameter σ goes to zero in accordance with the sample size. When the parameter σ diverges, generalization bounds can be also derived as shown below.

Theorem 15 *Suppose that Assumption 2 holds, $f^* \in \mathcal{H}$, and the risk functional $\mathcal{R}^\sigma(\cdot)$ is defined in association with a calibrated modal regression kernel K_σ and the corresponding representing function ϕ . Let $f_{\mathbf{z},\sigma}$ be produced by (2.4) with $\sigma > 1$. Then for any $0 < \delta < 1$, with probability at least $1 - \delta$, it holds that*

$$\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma}) \lesssim \frac{\log \delta^{-1}}{\sigma \sqrt{n}}.$$

Proof Similar to the proof of Theorem 14, the desired bound can be established by applying Lemma 13 to the random variable ξ in (3.2) with the only difference in bounding $\mathbb{E}\xi^2$. Recall that for a calibrated modal regression kernel K_σ , its representing function ϕ is bounded. Therefore, we have

$$\mathbb{E}\xi^2 = \mathbb{E} \left[\frac{1}{\sigma} \phi \left(\frac{Y - f_{\mathcal{H},\sigma}(X)}{\sigma} \right) - \frac{1}{\sigma} \phi \left(\frac{Y - f(X)}{\sigma} \right) \right]^2 \lesssim \sigma^{-2}.$$

In order to accomplish the proof, it suffices to apply Lemma 13 to the random variable ξ with $B = 2\|\phi\|_\infty/\sigma$, $\gamma = 0$, $c_p = \sigma^{-2p}$, and $c_\gamma = \sigma^{-2}$. By following the same procedure, the desired conclusion in Theorem 15 can be obtained. \blacksquare

The ERM learning scheme (2.4) is adaptive in that the scale parameter σ may vary in correspondence to the sample size n , e.g., $\sigma = n^\theta$ with $\theta \in \mathbb{R}$. Note from Theorem 15 that, with a properly chosen σ value, the ERM scheme (2.4) is generalization consistent in the sense that the generalization risk $\mathcal{R}^\sigma(f_{\mathbf{z},\sigma})$ converges to $\mathcal{R}^\sigma(f^*)$ when the sample size n tends to infinity. It is also interesting to note that a wide range of σ values is admitted to ensure such a consistency property as shown in the following corollary.

Corollary 16 *Suppose that Assumption 2 holds, $f^* \in \mathcal{H}$, and the risk functional $\mathcal{R}^\sigma(\cdot)$ is defined in association with a calibrated modal regression kernel K_σ and the representing function ϕ . Let $f_{\mathbf{z},\sigma}$ be produced by (2.4). Then for any $0 < \delta < 1$, with probability at least $1 - \delta$, it holds that*

$$\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma}) \rightarrow 0,$$

when $n \rightarrow +\infty$ and $\sigma := n^\theta$ with $\theta \in \left(-\frac{2}{2+3p}, +\infty\right)$.

Corollary 16 is an immediate result of Theorems 14 and 15 and its proof is omitted here. With a properly chosen σ value, the following conclusion reveals that the ERM scheme (2.4) is also modal regression consistent. This gives an affirmative answer to Question 2 listed in Subsection 3.1.

Theorem 17 *Suppose that Assumptions 1, 2 hold, and $f^* \in \mathcal{H}$. Let $f_{\mathbf{z},\sigma}$ be produced by (2.4) which is induced by a calibrated modal regression kernel K_σ with $\sigma = \mathcal{O}(n^{-\frac{2}{10+3p}})$. For any $0 < \delta < 1$, with probability at least $1 - \delta$, it holds that*

$$\mathcal{R}(f_{\mathbf{z},\sigma}) - \mathcal{R}(f^*) \lesssim n^{-\frac{4}{10+3p}} \log(\delta^{-1}).$$

Proof Since Assumption 2 holds, $f^* \in \mathcal{H}$, and K_σ is a calibrated modal regression kernel, from Theorem 14 we know that for any $0 < \delta < 1$, with probability at least $1 - \delta$, we have

$$\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma}) \lesssim \left(\frac{1}{n\sigma} + \frac{\sigma^{-\frac{2+3p}{4}}}{n^{1/2}} + \frac{\sigma^{-\frac{2+3p}{2+p}}}{n^{\frac{2}{2+p}}} \right) \log \left(\frac{1}{\delta} \right).$$

When Assumption 1 holds and K_σ is a calibrated modal regression kernel, Theorem 9 yields

$$\left| \{\mathcal{R}(f^*) - \mathcal{R}(f_{\mathbf{z},\sigma})\} - \{\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma})\} \right| \lesssim \sigma^2.$$

As a result, for any $0 < \delta < 1$, with probability at least $1 - \delta$, we have

$$\mathcal{R}(f^*) - \mathcal{R}(f_{\mathbf{z},\sigma}) \lesssim \sigma^2 + \left(\frac{1}{n\sigma} + \frac{\sigma^{-\frac{2+3p}{4}}}{n^{1/2}} + \frac{\sigma^{-\frac{2+3p}{2+p}}}{n^{\frac{2}{2+p}}} \right) \log \left(\frac{1}{\delta} \right).$$

With the choice $\sigma = \mathcal{O}(n^{-\frac{2}{10+3p}})$, the proof of Theorem 17 can be accomplished. \blacksquare

3.4. Towards the Function Estimation Calibration Problem

We now explore the relation between the modal regression consistency of $f_{\mathbf{z},\sigma}$ and its estimation consistency, which is termed as function estimation calibration problem in our study. From the studies in Heinrich (2013); Dearborn and Frongillo (2018), we realized that without further distributional assumptions, it is in general hopeless to learn the conditional mode through ERM approaches. In our study, we need to impose some further assumptions on the conditional density $p_{\epsilon|X}$ (see e.g., Doss and Wellner (2016)).

Definition 18 (Strongly s -Concave Density) *A density p is **strongly s -concave** if it exhibits one of the following forms:*

1. $p = \varphi_+^{1/s}$ for some strongly concave function φ if $s > 0$, where $\varphi_+ = \max\{\varphi, 0\}$;
2. $p = \exp(\varphi)$ for some strongly concave function φ if $s = 0$;
3. $p = \varphi_+^{1/s}$ for some strongly convex function φ if $s < 0$.

Assumption 3 *The density of ϵ conditioned on \mathcal{X} , denoted by $p_{\epsilon|X}(\cdot|X)$, satisfies the following conditions:*

1. $\sup_{x \in \mathcal{X}} p_{\epsilon|X}(0|X = x) = c_3$;
2. $p_{\epsilon|X}(t | X = x) \leq p_{\epsilon|X}(0|X = x)$, $\forall t \in \mathbb{R}, x \in \mathcal{X}$;
3. $\inf_{t \in [-2M, 2M], x \in \mathcal{X}} p_{\epsilon|X}(t|X = x) = c_0 > 0$;
4. $p_{\epsilon|X}(\cdot | X)$ denotes strongly s -concave densities for all realizations of X .

Conditions 1 and 2 in Assumption 3 require that the global mode of the conditional density $p_{\epsilon|X}$ for any realization of X in \mathcal{X} is uniquely zero while Condition 3 rules out densities that are not bounded away from below in the vicinity of this unique mode. The first two conditions hold for continuous densities with a unique global mode. Condition 4 assumes the strongly s -concave density assumption on $p_{\epsilon|X}$, which is typical from a statistical viewpoint as it holds for common symmetric and skewed distributions. Several representative examples are listed below:

Example 1 (Student's t -distribution) *Let ρ be a Student's t -distribution. Its probability density function p is*

$$p(t) = \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}},$$

where ν is the number of degrees of freedom and Γ is the gamma function. Specifically, when $\nu = 1$, it gives the density function of a typical heavy-tailed distribution, namely, Cauchy distribution; when $\nu = \infty$, it is the density function of a most common probability distribution, i.e., Gauss distribution. One can easily see that for Student's t -distributions, their densities are strongly s -concave and are of the form 3 in Definition 18.

Example 2 (Skewed normal distribution) *Let ρ be a skewed normal distribution with the probability density function*

$$p(t|\mu, \theta, \tau) = \frac{4\tau(1-\tau)}{\sqrt{2\pi}\theta^2} \exp\left\{-\frac{2(x-\mu)^2}{\sigma^2} (\tau - \mathbb{1}_{(x \leq \mu)}(x))\right\},$$

where $\mathbb{1}_A(x)$ is the indicator function that takes the value 1 if A is true and 0, otherwise. Clearly, the above density is also strongly s -concave and is of the form 2 in Definition 18.

When Assumption 3 holds, the function estimation convergence can be elicited from the convergence of the modal regression risk, as shown in the following theorem.

Theorem 19 *Suppose that Assumption 3 holds and let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be a measurable function in \mathcal{H} . Then, it holds that*

$$\|f - f^*\|_{L^2_{\rho_X}}^2 \lesssim \mathcal{R}(f^*) - \mathcal{R}(f).$$

Proof If Assumption 3 is fulfilled, then $p_{\epsilon|X}$ is strongly s -concave. We verify the desired relation by discussing different cases of s . If $s = 0$, we know that $-\log p_{\epsilon|X}$ is strongly convex for all x . Consequently, in this case, it holds that

$$\begin{aligned} \|f - f^*\|_{L^2_{\rho_X}}^2 &\lesssim \int_{\mathcal{X}} [-\log p_{\epsilon|X}(f(x) - f^*(x) | X = x) + \log p_{\epsilon|X}(0 | X = x)] d\rho_X(x) \\ &\lesssim \int_{\mathcal{X}} [p_{\epsilon|X}(0 | X = x) - p_{\epsilon|X}(f(x) - f^*(x) | X = x)] d\rho_X(x), \end{aligned}$$

where the last inequality is a consequence of the mean value theorem and Assumption 3. If $s > 0$, $-p_{\epsilon|X}^s$ is strongly convex for all x , then

$$\begin{aligned} \|f - f^*\|_{L_{\rho_X}^2}^2 &\lesssim \int_{\mathcal{X}} [-p_{\epsilon|X}^s(f(x) - f^*(x) | X = x) + p_{\epsilon|X}^s(0 | X = x)] d\rho_X(x) \\ &\lesssim \max\{sc_0^{s-1}, sc_3^{s-1}\} \int_{\mathcal{X}} [p_{\epsilon|X}(0 | X = x) - p_{\epsilon|X}(f(x) - f^*(x) | X = x)] d\rho_X(x), \end{aligned}$$

where the second inequality is due to the Lipschitz continuity of $h(t) = t^s$ and Assumption 3. If $s < 0$, $p_{\epsilon|X}^s$ is strongly convex for all x . In this case, we have

$$\begin{aligned} \|f - f^*\|_{L_{\rho_X}^2}^2 &\lesssim \int_{\mathcal{X}} [p_{\epsilon|X}^s(f(x) - f^*(x) | X = x) - p_{\epsilon|X}^s(0 | X = x)] d\rho_X(x) \\ &\lesssim -sc_0^{s-1} \int_{\mathcal{X}} [p_{\epsilon|X}(0 | X = x) - p_{\epsilon|X}(f(x) - f^*(x) | X = x)] d\rho_X(x), \end{aligned}$$

where the second inequality is again due to the Lipschitz continuity of $h(t) = t^s$ and Assumption 3. Recalling the fact that

$$\mathcal{R}(f^*) - \mathcal{R}(f) = \int_{\mathcal{X}} [p_{\epsilon|X}(0 | X = x) - p_{\epsilon|X}(f(x) - f^*(x) | X = x)] d\rho_X(x),$$

we complete the proof of Theorem 19. ■

Combining the estimates established in the above several subsections, we are now able to answer Question 3 raised in Subsection 3.1.

Theorem 20 *Suppose that Assumptions 1, 2, and 3 hold, and $f^* \in \mathcal{H}$. Let $f_{z,\sigma}$ be produced by (2.4) which is induced by a calibrated modal regression kernel K_σ with $\sigma = \mathcal{O}(n^{-\frac{2}{10+3p}})$. For any $0 < \delta < 1$, with probability at least $1 - \delta$, we have*

$$\|f_{z,\sigma} - f^*\|_{L_{\rho_X}^2}^2 \lesssim n^{-\frac{4}{10+3p}} \log(\delta^{-1}).$$

Proof The theorem can be proved by combining the estimates in Theorems 17 and 19. ■

3.5. Some Remarks

We give some remarks here. As noted earlier, most of the existing studies on modal regression were conducted by resorting to maximizing the joint density estimator or the conditional density estimator. However, there are two main barriers when seeking the maximizer in this way. First, from a statistical learning viewpoint, learning the maximizer of the joint density or the conditional density is a local type learning scheme, in which one has to train the model for each test point. Second, the estimation of a high-dimensional joint or conditional density may suffer from the curse of dimensionality. In our proposed ERM approach to modal regression, the hypothesis space \mathcal{H} is a function space that can be

infinite-dimensional. In practice, it can be specified by applying certain regularization procedures. Moreover, the prevalent kernel-based methods can be naturally integrated since the hypothesis space can be chosen as a subset of a certain reproducing kernel Hilbert space. On the other hand, the proposed ERM approach to modal regression only involves a one-dimensional density estimation problem. From the above comparisons and the learning theory analysis conducted in this paper, it is easy to see that our study provides a different take on modal regression and the proposed ERM approach distinguishes our work with the existing studies.

4. Modal Regression Interpretation of Correntropy based Regression

As mentioned above, our study on modal regression in this paper is initiated to understand the so-called maximum correntropy criterion in regression problems (see Liu et al., 2007; Principe, 2010). In this sense, the present study is a continuation of our previous work in Feng et al. (2015). As a generalized correlation measurement, correntropy has been drawing much attention recently. Owing to its prominent merits on robustness, it has been pervasively used and has found many real-world applications in signal processing, machine learning, and computer vision (see e.g., Bessa et al., 2009; He et al., 2011, 2012; Lu et al., 2013; Chen et al., 2016a).

4.1. Correntropy and Correntropy based Regression

Mathematically speaking, correntropy is a generalized similarity measure between two scalar random variables U and V , which is defined by $\mathcal{R}^\sigma(U, V) = \mathbb{E}K_\sigma(U, V)$. Here K_σ is a Gaussian kernel given by $K_\sigma(u, v) = \exp\{-(u - v)^2/\sigma^2\}$ with the bandwidth $\sigma > 0$, (u, v) being a realization of (U, V) . Given a set of i.i.d observations $\mathbf{z} = \{(x_i, y_i)\}_{i=1}^n$, for any $f : \mathcal{X} \rightarrow \mathbb{R}$, the empirical estimator of the correntropy between $f(X)$ and Y is given as

$$\mathcal{R}_n^\sigma(f) := \frac{1}{n} \sum_{i=1}^n K_\sigma(y_i, f(x_i)).$$

The maximum correntropy criterion based regression models the empirical target function by maximizing the empirical estimator of the correntropy \mathcal{R}^σ as follows

$$f_{\mathbf{z}, \sigma} = \arg \max_{f \in \mathcal{H}} \mathcal{R}_n^\sigma(f), \tag{4.1}$$

where \mathcal{H} is assumed to be a compact subset of $C(\mathcal{X})$. Here, $C(\mathcal{X})$ is denoted as the Banach space of continuous functions on \mathcal{X} . The maximum correntropy criterion in regression problems has shown its efficiency for cases where non-Gaussian noise or outliers are present (see e.g., Liu et al., 2007; Principe, 2010; Wang et al., 2013).

In the literature, existing understanding of the maximum correntropy criterion and MCCR is still limited. More frequently, the maximum correntropy criterion is roughly taken as a robustified least squares criterion, analogously to the trimmed least squares criterion. However, the statistical performance of $f_{\mathbf{z}, \sigma}$ and its relation to the least squares criterion are not clear. The barriers are mainly caused by the presence of the scale parameter σ and the non-convexity of the related model. Recently, some theoretical understanding towards

the maximum correntropy criterion was conducted in Feng et al. (2015) by introducing a distance-based regression loss, the study of which is inspired by those on information theoretic learning in Hu et al. (2013) and Fan et al. (2016). The main conclusion drawn in Feng et al. (2015) is that MCCR is essentially robustified mean regression with diverging σ values. On the other hand, our study conducted in this paper shows that with diminishing $\sigma(n)$ values, MCCR is, in fact, modal regression. The built-in robustness of modal regression schemes may explain the empirical successes of MCCR from a different viewpoint.

4.2. A General Picture of Correntropy based Regression

Based on this study and the study in Feng et al. (2015), we are now able to depict a general picture of the correntropy based regression from a statistical learning viewpoint. To this end, we expost the correntropy based regression by considering three different cases below, namely, (1): $\sigma = \sigma(n) \rightarrow \infty$; (2): $\sigma := \sigma_0$ for some $\sigma_0 > 0$, that is, σ is fixed and independent of the sample size n ; (3): $\sigma := \sigma(n) \rightarrow 0$. Before proceeding, we recall the following data-generating model

$$Y = f^*(X) + \epsilon.$$

We first consider the case when $\sigma(n) \rightarrow \infty$. Under the zero-mean noise assumption on ϵ , i.e., $\mathbb{E}(\epsilon|X) = 0$, MCCR (4.1) with $\sigma(n) \rightarrow \infty$ encourages the approximation of $f_{\mathbf{z},\sigma}$ towards the conditional mean function $\mathbb{E}(Y|X)$ and the scale parameter σ in this case plays a trade-off role between robustness and generalization. More explicitly, in this case, MCCR is mean regression calibrated in the sense of the following theorem, see also Lemma 7 in Feng et al. (2015):

Theorem 21 (Lemma 7, Feng et al. (2015)) *Assume that $\mathbb{E}Y^4 < \infty$ and denote $f^* = \mathbb{E}(Y|X)$. For any $f \in \mathcal{H}$, it holds that*

$$\left| \|f - f^*\|_{L^2_{\rho_X}}^2 - |\sigma^3(\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f))| \right| \lesssim \sigma^{-2}.$$

It turns out that when $\sigma(n)$ is properly chosen with $\sigma(n) \rightarrow \infty$, the consistency of $\mathcal{R}^\sigma(f^*) - \mathcal{R}^\sigma(f_{\mathbf{z},\sigma})$ implies the consistency of $\|f_{\mathbf{z},\sigma} - f^*\|_{L^2_{\rho_X}}^2$. Moreover, the following convergence rates are established in Feng et al. (2015):

Theorem 22 *Assume that $f^* = \mathbb{E}(Y|X) \in \mathcal{H}$ and $\mathbb{E}Y^4 < +\infty$. Under a mild capacity assumption on \mathcal{H} , for any $0 < \delta < 1$, with confidence $1 - \delta$, it holds that*

$$\|f_{\mathbf{z},\sigma} - f^*\|_{L^2_{\rho_X}}^2 \lesssim \sigma^{-2} + \sigma n^{-1/(1+p)},$$

where the index $p > 0$ reflects the capacity of the hypothesis space \mathcal{H} .

Obviously, according to the above theorem, when σ is chosen as $\sigma := n^{-1/(3+3p)}$, the convergence rates for $\|f_{\mathbf{z},\sigma} - f^*\|_{L^2_{\rho_X}}^2$ of the type $\mathcal{O}(n^{-2/(3+3p)})$ can be established. It is worth to mention that recently, in Feng and Wu (2019), the above moment condition $\mathbb{E}Y^4 < +\infty$ was further relaxed to $\mathbb{E}|Y|^{1+\zeta} < +\infty$ with $\zeta > 0$. Notice that in this case the underlying truth f^* corresponds to the conditional mean. Therefore, MCCR in this case

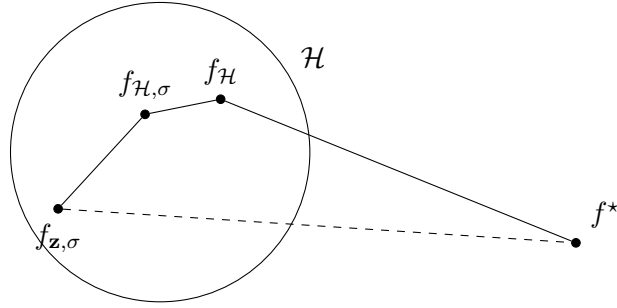


Figure 2: A schematic illustration of the mechanism of correntropy-based regression when $\sigma(n) \rightarrow \infty$ and the noise variable ϵ is assumed to be zero-mean. $f_{\mathcal{H},\sigma}$ is the data-free counterpart of $f_{\mathbf{z},\sigma}$, $f_{\mathcal{H}}$ is the data-free least squares regression estimator and f^* is the conditional mean function $\mathbb{E}(Y|X)$.

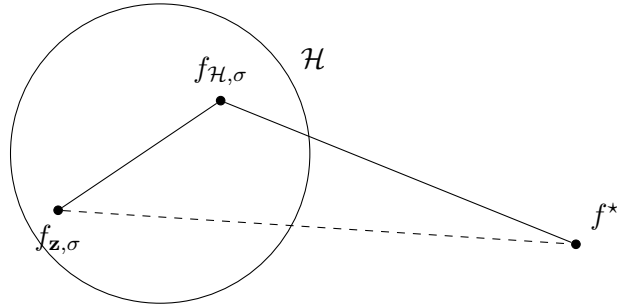


Figure 3: A schematic illustration of the mechanism of correntropy-based regression when σ is fixed and independent on n and the noise variable ϵ is assumed to be zero-mean. $f_{\mathcal{H},\sigma}$ is the data-free counterpart of $f_{\mathbf{z},\sigma}$ and f^* is the conditional mean function $\mathbb{E}(Y|X)$ or the conditional median function $\text{median}(Y|X)$.

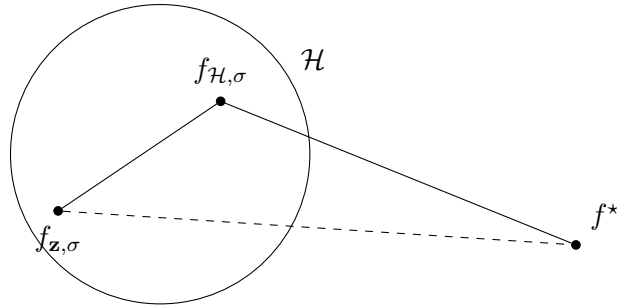


Figure 4: A schematic illustration of the mechanism of correntropy-based regression when $\sigma(n) \rightarrow 0$ and the noise variable ϵ is assumed to admit a unique global zero-mode. $f_{\mathcal{H},\sigma}$ is the data-free counterpart of $f_{\mathbf{z},\sigma}$ and f^* is the conditional mode function $\text{mode}(Y|X)$.

is essentially robustified mean regression. A schematic illustration of MCCR in this case is given in Fig. 2, in which $f_{\mathcal{H},\sigma}$ is the population version of $f_{\mathbf{z},\sigma}$ and $f_{\mathcal{H}}$ is the data-free least squares regression estimator. As argued in Feng et al. (2015), compared with the least squares regression, an additional bias, i.e., the distance between $f_{\mathcal{H},\sigma}$ and $f_{\mathcal{H}}$, appears when bounding the $L^2_{\rho_{\mathcal{X}}}$ -distance between $f_{\mathbf{z},\sigma}$ and the conditional mean function $\mathbb{E}(Y|X)$. Moreover, this bias in some sense reflects the trade-off between the convergence rate of $\|f_{\mathbf{z},\sigma} - f^*\|_{L^2_{\rho_{\mathcal{X}}}}^2$ and the robustness of $f_{\mathbf{z},\sigma}$. These observations were further justified in a regularized learning setup in Lv and Fan (2019).

	$\sigma(n) \rightarrow \infty$	σ fixed	$\sigma(n) \rightarrow 0$
resulting estimator	conditional mean estimator	conditional mean or median estimator	conditional mode estimator
target function	$\mathbb{E}(Y X)$	$\mathbb{E}(Y X)$ or $\text{median}(Y X)$	$\text{mode}(Y X)$
noise condition	weak moment condition	bounded symmetric or symmetric stable	allow skewness or heavy-tailedness
rates	$\mathcal{O}(n^{-2/(3+3p)})$	$\mathcal{O}(n^{-2/(2+p)})$	$\mathcal{O}(n^{-4/(10+3p)})$

Table 2: An overview of the three scenarios in correntropy based regression

The case when $\sigma = \sigma_0$, i.e., σ is fixed and independent of n , was investigated in Feng et al. (2015), Feng and Wu (2019), and Feng and Ying (2019). As argued in Feng and Wu (2019), with a fixed parameter σ and without imposing any noise assumptions, it is impossible to learn the truth function f^* . It turns out that in this case, if some noise assumptions are introduced, correntropy based regression regresses towards the conditional mean or the conditional median. More specifically, according to Lemma 18 in Feng et al. (2015), under bounded symmetric noise assumptions, it is also calibrated mean regression when σ_0 is properly chosen. Convergence rates of $\|f_{\mathbf{z},\sigma} - f^*\|_{L^2_{\rho_{\mathcal{X}}}}^2$ can be also established under such noise assumptions, see Theorem 6 in Feng et al. (2015). Inspired by the work in Fan et al. (2016), it is demonstrated in Feng and Ying (2019) that under the symmetric stable noise assumption, correntropy based regression can learn the underlying truth function f^* well where the truth function in this scenario corresponds to the conditional mean or the conditional median function.

The fact that MCCR can be cast as a modal regression problem when $\sigma(n) \rightarrow 0$ switches our attention from robust mean regression in Feng et al. (2015) to modal regression in this study. To recap, the modal regression scheme (2.4) with the Gaussian kernel as the modal regression kernel retrieves MCCR (4.1). From the arguments in the preceding sections, we know that under the assumption that the noise variable admits a unique global zero-mode, MCCR (4.1) with $\sigma(n) \rightarrow 0$ is modal regression calibrated. That is, under proper assumptions as listed in Theorem 20, one may expect the learning theory type convergence from the MCCR estimator to the modal regression function $\text{mode}(Y|X)$. Results reported in the above sections reveal that the modal regression problem can be also studied from an empirical risk minimization viewpoint. A schematic illustration of the mechanism of

correntropy-based regression when $\sigma(n) \rightarrow 0$ is presented in Fig. 4. In this case, the robustness of MCCR stems from the built-in robustness of modal regression estimators.

An overview of the above-discussed three scenarios in correntropy based regression is summarized in Table 2. To sum up, in short, what makes MCCR so special is that it results an interesting walk between modal regression and robustified mean regression by adjusting the scale parameter σ in correspondence to the sample size n .

5. Model Selection and Numerical Validations

This section is concerned with the implementation issues of the proposed ERM approach to modal regression. The model selection problem will be tackled by tailoring the technique of cross validation. Numerical validations on the effectiveness of the proposed modal regression estimators will also be provided.

5.1. Experimental Setup

In our empirical studies, the hypothesis space \mathcal{H} is chosen as a bounded subset of a reproducing kernel Hilbert space $\mathcal{H}_{\mathcal{K}}$ that is induced by a Mercer kernel \mathcal{K} . Specifically, we employ the following Tikhonov regularization to determine the radius of the working hypothesis space automatically:

$$f_{\mathbf{z},\sigma} := \arg \min_{f \in \mathcal{H}_{\mathcal{K}} \oplus \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \ell_{\sigma}(y_i - f(x_i)) + \lambda \|f\|_{\mathcal{K}}^2, \quad (5.1)$$

where ℓ_{σ} is the loss function $\ell_{\sigma}(t) = \sigma^2(1 - \exp(-t^2/\sigma^2))$, and $\lambda > 0$ is a regularization parameter. The representer theorem ensures that $f_{\mathbf{z},\sigma}$ can be modeled by

$$f_{\mathbf{z},\sigma}(x) = \sum_{i=1}^n \alpha_{\mathbf{z},i} \mathcal{K}(x, x_i) + b_{\mathbf{z}}, \quad x \in \mathbb{R},$$

where $\alpha_{\mathbf{z}} = (\alpha_{\mathbf{z},1}, \dots, \alpha_{\mathbf{z},n})^{\top} \in \mathbb{R}^n$ and $b_{\mathbf{z}} \in \mathbb{R}$ are learned from (5.1). For the Mercer kernel \mathcal{K} , we use the Gaussian kernel $\mathcal{K}(x, x') = \exp(-\|x - x'\|^2/h^2)$ with the bandwidth parameter $h > 0$.

5.2. Algorithms

The regularization problem (5.1) is essentially a regularized M-estimation problem. We, therefore, apply the iteratively re-weighted least squares algorithm to solve it. The pseudocode of the iteratively re-weighted least squares algorithm is listed in Algorithm 1. For each iteration in Algorithm 1, the weight is updated as follows:

$$\omega_i^{k+1} = \frac{|\nabla \ell_{\sigma}(y_i - \mathcal{K}_i^{\top} \alpha^k - b^k)|}{|y_i - \mathcal{K}_i^{\top} \alpha^k - b^k|}, \quad i = 1, \dots, n, \quad (5.2)$$

with the initial guess α^0, b^0 being zero.

Algorithm 1: Iteratively Re-weighted Least Squares Algorithm for Solving (5.1)

Input: data $\{(x_i, y_i)\}_{i=1}^n$, regularization parameter $\lambda > 0$, Gaussian kernel bandwidth $h > 0$, scale parameter $\sigma > 0$ and the initial guess $\alpha^0 \in \mathbb{R}^n$, $b^0 \in \mathbb{R}$.

Output: the learned coefficient $\alpha^{k+1} = (\alpha_1^{k+1}, \dots, \alpha_n^{k+1})^\top$ and $b^{k+1} \in \mathbb{R}$.

while the stopping criterion is not satisfied **do**

- Compute α^{k+1} and b^{k+1} by solving the following weighted least squares problem:

$$(\alpha^{k+1}, b^{k+1}) = \arg \min_{\alpha \in \mathbb{R}^n, b \in \mathbb{R}} \sum_{i=1}^n \omega_i^{k+1} (y_i - \mathcal{K}_i^\top \alpha - b)^2 + \lambda \alpha^\top \mathcal{K} \alpha,$$

where ω_i^{k+1} is specified in (5.2).

- Set $k := k + 1$.

end while

5.3. Model Selection via Concatenated Cross Validation

We now discuss the model selection problem of the proposed modal regression estimator. Here, the problem of model selection refers to the selection of the three tuning parameters, i.e., the regularization parameter λ , the bandwidth parameter h of the Gaussian kernel, and the scale parameter σ in the loss function.

In our study, we choose these parameters by tailoring the frequently used cross-validation technique and propose Concatenated Cross Validation (CCV) for model selection. In order to carry out the cross-validation process, we need to choose an error criterion. As we are interested in learning the conditional mode function, the mean squared error criterion, the absolute deviation error criterion, as well as the criteria under robustness constraints, see e.g., Cantoni and Ronchetti (2001), may not serve well for this purpose. Recall that the ERM approach for modal regression we proposed in this study can be also re-expressed as follows

$$f_{\mathbf{z}, \sigma} = \arg \max_{f \in \mathcal{H}} \frac{1}{n\sigma} \sum_{i=1}^n \exp \left(-\frac{(y_i - f(x_i))^2}{\sigma^2} \right),$$

where the hypothesis space \mathcal{H} is chosen as a subset of a reproducing kernel Hilbert space induced by the Gaussian kernel as mentioned above. The criterion that we use in CCV is essentially the loss function in the above ERM scheme. More explicitly, denoting $\{(x_i, y_i)\}_{i=1}^m$ as the validation set and $\{\hat{y}_{i, \sigma}\}_{i=1}^m$ the estimated values, CCV can be proceeded through the following steps:

Step 1: We implement a first five-fold cross validation under the following criterion

$$\arg \max_{\sigma} \frac{1}{m\sigma_0} \sum_{i=1}^m \exp \left(-\frac{(y_i - \hat{y}_{i, \sigma})^2}{\sigma_0^2} \right),$$

where the initial value σ_0 is set as $m^{-1/5}$, which is the optimal σ value according to our theoretical analysis. We denote the best σ value selected in this step as σ_1 .

Step 2: We then implement a second five-fold cross validation under the following updated criterion

$$\arg \max_{\sigma} \frac{1}{m\sigma_1} \sum_{i=1}^m \exp \left(-\frac{(y_i - \hat{y}_{i,\sigma})^2}{\sigma_1^2} \right).$$

We denote the best σ value selected in this step as σ_2 .

Step 3: We continue to implement a third five-fold cross validation under the following updated criterion

$$\arg \max_{\sigma} \frac{1}{m\sigma_2} \sum_{i=1}^m \exp \left(-\frac{(y_i - \hat{y}_{i,\sigma})^2}{\sigma_2^2} \right).$$

We denote the best σ value selected in this step as σ_3 . Note that in the above steps, the estimated values $\{\hat{y}_{i,\sigma}\}_{i=1}^m$ also depend on the tuning parameters λ and h , which are also updated accordingly at each step. We suppress the two subscripts for simplification.

Step 4: With the selected σ value in Step 3, we then train the regularized ERM model by using the iterative reweighted least squares algorithm. We then take the resulting estimator as the modal regression estimator and proceed with the prediction process.

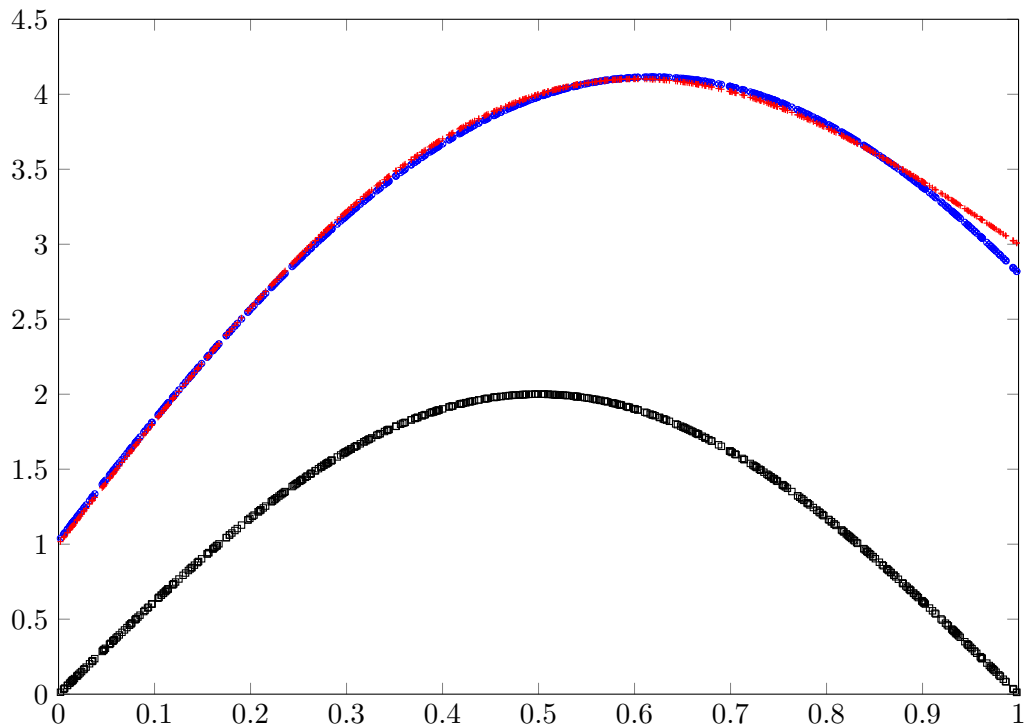


Figure 5: The dotted red curve with square marks is the conditional mode function f_{MO} for observations generated by (5.3) while the dotted black curve with plus marks gives the conditional mean function f_{ME} . The dotted blue curve with \otimes marks represents the learned estimator $f_{z,\sigma}$ from noisy observations.

5.4. Numerical Validation on a Toy Example

We validate the effectiveness of the proposed modal regression estimator on the following toy example. We generate artificial data through the following regression model

$$y = f^*(x) + \kappa(x)\epsilon, \quad (5.3)$$

where $x \sim U(0, 1)$, $f^*(x) = 2 \sin(\pi x)$, and $\kappa(x) = 1 + 2x$. The noise variable is distributed as $\epsilon \sim 0.5N(-1, 4^2) + 0.5N(1, 0.1^2)$. A similar example was employed in Yao and Li (2014). With simple calculations, it is easy to see that the conditional mean function is $f_{\text{ME}} = 2 \sin(\pi x)$ and the conditional mode function is approximately $f_{\text{MO}} = 2 \sin(\pi x) + 1 + 2x$. In our experiment, 600 observations are drawn from the above data-generating model and the size of the test set is also set to 600. The reconstructed curve is plotted at the test points in Fig. 5, in which the conditional mean function f_{ME} and the conditional mode function f_{MO} are also plotted for comparisons. In our experiment, we choose the three tuning parameters, i.e., the bandwidth parameter h of the Gaussian kernel, the regularization parameter λ , and the scale parameter σ in the loss function, by using Concatenated Cross Valuation described above.

From Fig. 5, it is easy to see that the proposed modal regression estimator $f_{\mathbf{z}, \sigma}$ can learn the conditional mode function f_{MO} well instead of learning the conditional mean function f_{ME} . It is interesting to point out that the obtained empirical target function $f_{\mathbf{z}, \sigma}$ can also learn the conditional mean function with a large σ value as explained in Section 4.

5.5. Application to Speed-Flow Data

We now apply the proposed modal regression estimator to speed-flow data. Speed-flow data are intensively discussed in transportation science, which are usually visualized in terms of speed-flow diagrams. In this subsection, we apply the proposed modal regression approach to the analysis of the speed-flow data collected in Petty et al. (1996), the speed-flow diagrams of which are presented in Figs. 6 and 7. In the speed-flow diagrams, the x -axis is traffic flow that is measured in vehicles per lane per hour while the y -axis is speed measured in miles per hour. The speed-flow data analyzed here contain two data sets collected in 1993 on two individual lanes (lane 2 and lane 3) of the 4-lane Californian freeway I-880. The data were collected by loop detectors, and the time units are 30 seconds per observation, see Einbeck and Tutz (2006) for more background details. This speed-flow data contains 1318 observations and are publicly available in the R-package `hdrcde`. From the speed-flow diagrams, it can be observed that the mean regression function may not be able to characterize the functional relation between speed and traffic flow. This is also observed in many related studies that analyze the speed-flow data, see e.g., Einbeck and Tutz (2006). This is because the less dense cloud of data points at the bottom of the two figures, which corresponds to situations where speed is dismissed, may be interpreted as abnormal observations when pursuing such a functional relation.

In our experiments, we apply the proposed modal regression approach to pursuing the functional relation. By following the same setup as in our above experiments on artificial data, we plot the learned modal regression estimator as well as the mean regression estimator resulting from kernel ridge regression. From the reported experimental results in Fig. 6 and Fig. 7, it can be seen that modal regression estimator is less sensitive to abnormal

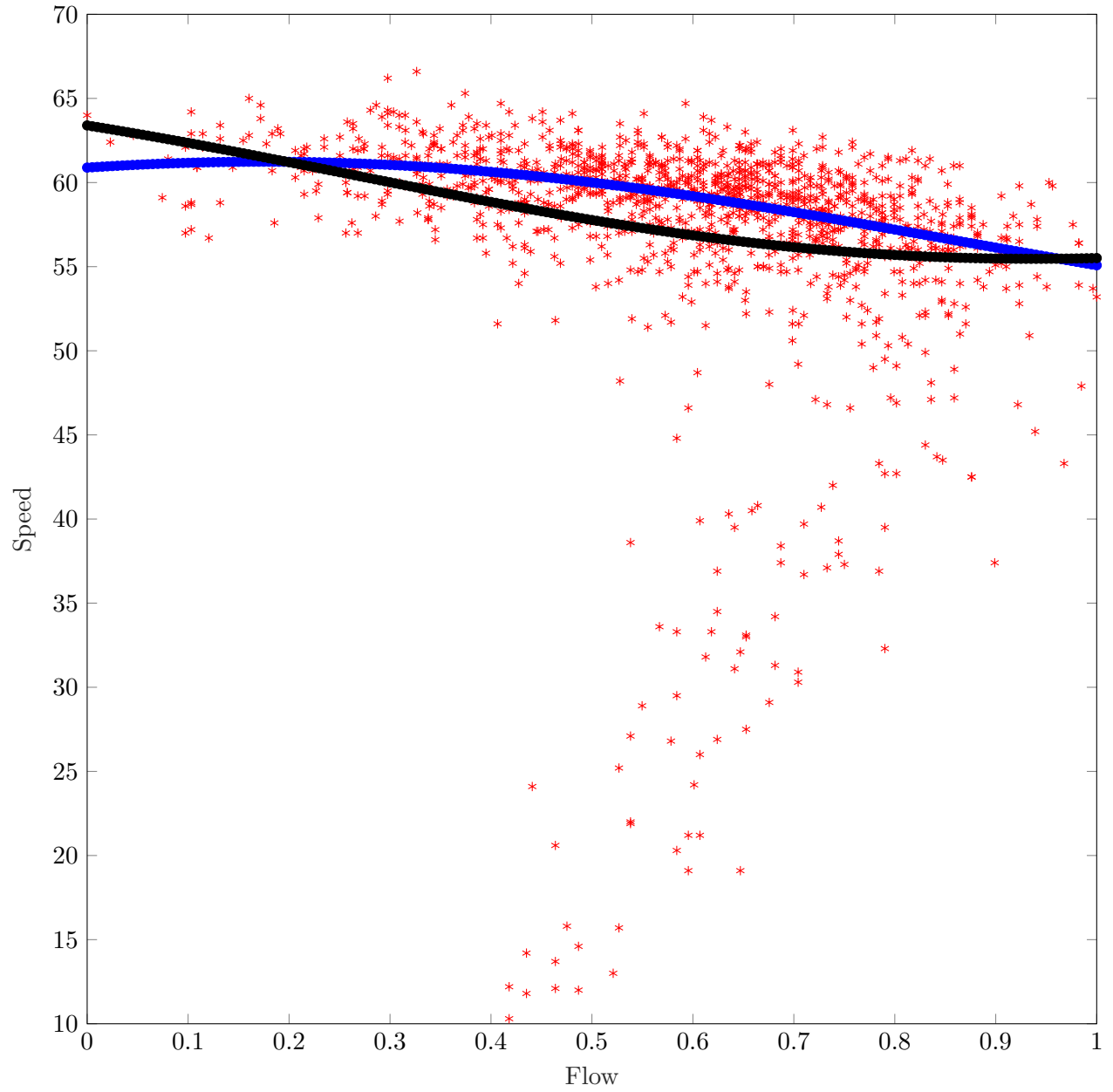


Figure 6: The blue curve represents the conditional mode function estimator $f_{z,\sigma}$ for 1318 observations of lane 2 while the black curve gives the conditional mean function estimator by kernel ridge regression.

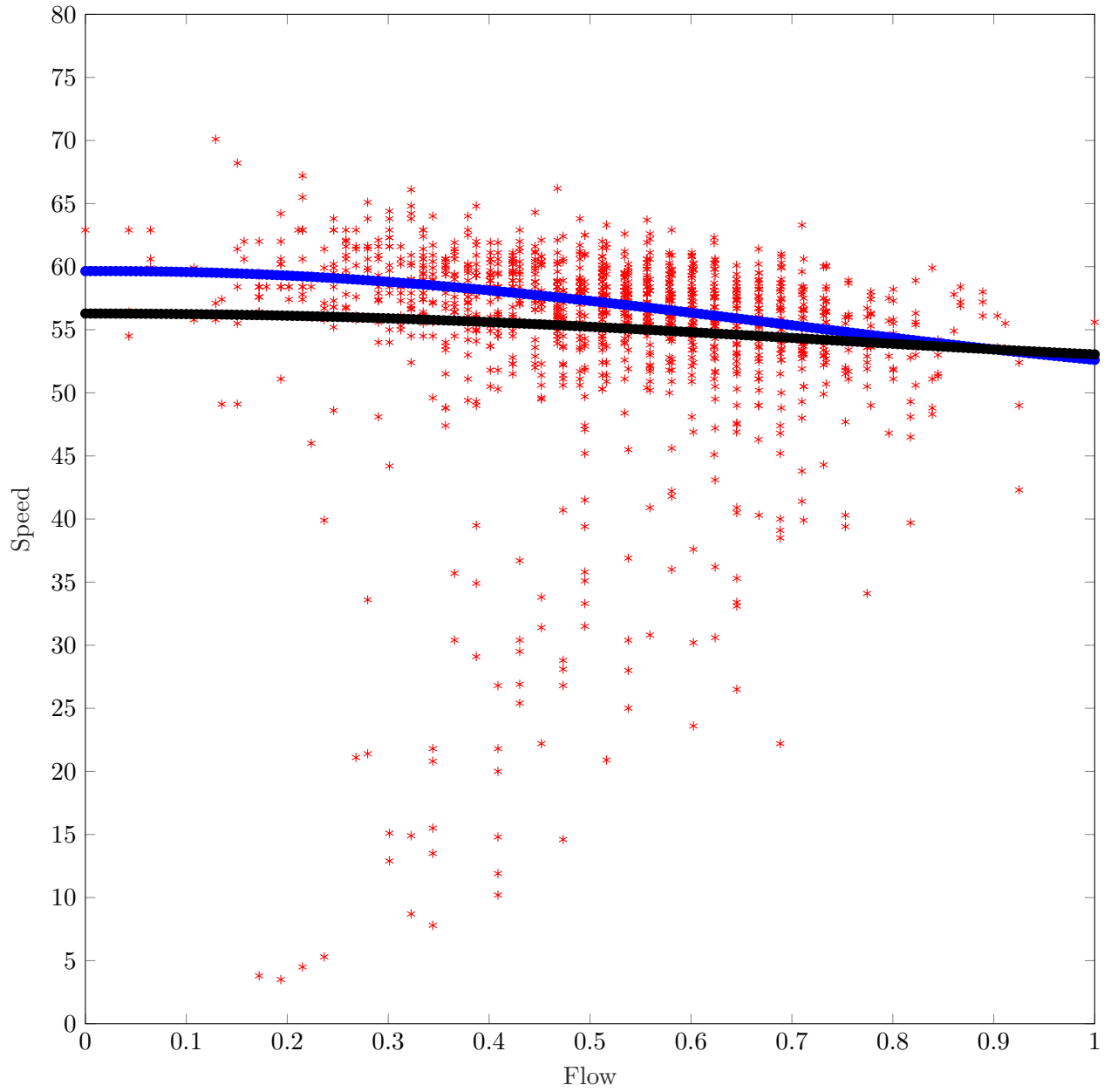


Figure 7: The blue curve represents the conditional mode function estimator $f_{z,\sigma}$ for 1318 observations of lane 3 while the black curve gives the conditional mean function estimator by kernel ridge regression.

observations and serves better in trend estimation when analyzing speed-flow data. It would be interesting to explore more real-world applications of the modal regression estimator learned through the proposed ERM approach, which will be the future work of our study in this respect.

6. Conclusions

As one of the important regression protocols, modal regression has not been much studied yet in the statistical learning literature. In this study, we investigated the modal regression problem from a statistical learning viewpoint. By assuming the existence and the uniqueness of the global mode of the conditional distribution in regression, we reformulated the modal regression problem into the classical empirical risk minimization framework. In particular, such a reformulation renders the associated modal regression approach dimension-independent. A learning theory framework for analyzing and assessing the proposed modal regression estimator was also developed. Based on the proposed statistical learning treatment on modal regression, we gained some insights into the regression problem. These insights include: first, modal regression problem can be tackled via empirical risk minimization and can be also interpreted from a kernel density estimation point of view; second, learning for modal regression is generalization consistent and modal regression calibrated in the sense defined in our study; third, function estimation consistency and convergence in the sense of the $L^2_{\rho_X}$ -distance can be derived in modal regression. These findings in return unveil the working mechanism of MCCR when its scale parameter tends to zero as in this case, it corresponds to a modal regression problem.

Acknowledgments

The authors would like to thank the Action Editor and the reviewers for their constructive suggestions and comments that improved the quality of this paper. The research leading to these results has received funding from the European Research Council under the European Union’s Seventh Framework Programme (FP7/2007-2013) / ERC AdG A-DATADRIVE-B (290923) and ERC AdG E-DUALITY (787960) under the European Union’s Horizon 2020 research and innovation programme. This paper reflects only the authors’ views, the Union is not liable for any use that may be made of the contained information. Research Council KUL: GOA/10/09 MaNet, CoE PFV/10/002 (OPTEC), BIL12/11T; PhD/Postdoc grants. Flemish Government: FWO: projects: G.0377.12 (Structured systems), G.088114N (Tensor based data similarity); PhD/Postdoc grants. IWT: projects: SBO POM (100031); PhD/Postdoc grants. iMinds Medical Information Technologies SBO 2014. Belgian Federal Science Policy Office: IUAP P7/19 (DYSCO, Dynamical systems, control and optimization, 2012-2017). Yunlong Feng also gratefully acknowledges the support of Simons Foundation Collaboration Grant #572064 and the Ralph E. Powe Junior Faculty Enhancement Award by Oak Ridge Associated Universities. The research of Jun Fan was supported in part by the Hong Kong RGC Early Career Schemes 22303518, and the NSF grant of China (No. 11801478). The corresponding author is Jun Fan.

References

- Martin Anthony and Peter L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, 2009.
- Markus Baldauf and Joao Santos Silva. On the use of robust regression in econometrics. *Economics Letters*, 114(1):124–127, 2012.
- Peter L. Bartlett, Olivier Bousquet, and Shahar Mendelson. Local Rademacher complexities. *The Annals of Statistics*, 33(4):1497–1537, 2005.
- Peter L. Bartlett, Michael I. Jordan, and Jon D. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- Ricardo J. Bessa, Vladimiro Miranda, and Joao Gama. Entropy and correntropy against minimum square error in offline and online three-day ahead wind power forecasting. *IEEE Transactions on Power Systems*, 24(4):1657–1666, 2009.
- Eva Cantoni and Elvezio Ronchetti. Resistant selection of the smoothing parameter for smoothing splines. *Statistics and Computing*, 11(2):141–146, 2001.
- Badong Chen, Lei Xing, Haiquan Zhao, Nanning Zheng, and José C. Principe. Generalized correntropy for robust adaptive filtering. *IEEE Transactions on Signal Processing*, 64(13):3376–3387, 2016a.
- Yen-Chi Chen, Christopher R. Genovese, Ryan J. Tibshirani, and Larry Wasserman. Non-parametric modal regression. *The Annals of Statistics*, 44(2):489–514, 2016b.
- Herman Chernoff. Estimation of the mode. *Annals of the Institute of Statistical Mathematics*, 16(1):31–41, 1964.
- Grard Collomb, Wolfgang Härdle, and Salima Hassani. A note on prediction via estimation of the conditional mode function. *Journal of Statistical Planning and Inference*, 15(2):227–236, 1987.
- Dorin Comaniciu and Peter Meer. Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(5):603–619, 2002.
- Felipe Cucker and Ding-Xuan Zhou. *Learning Theory: An Approximation Theory Viewpoint*. Cambridge University Press, 2007.
- Sanjoy Dasgupta and Samory Kpotufe. Optimal rates for k -NN density and mode estimation. In *Advances in Neural Information Processing Systems*, pages 2555–2563, 2014.
- Krisztina Dearborn and Rafael Frongillo. On the indirect elicibility of the mode and modal interval. *Annals of the Institute of Statistical Mathematics*, pages 1–14, 2018.
- Charles R. Doss and Jon A. Wellner. Global rates of convergence of the MLEs of log-concave and s -concave densities. *The Annals of Statistics*, 44(3):954–981, 2016.

- William F. Eddy. Optimum kernel estimators of the mode. *The Annals of Statistics*, 8(4): 870–882, 1980.
- Jochen Einbeck and Gerhard Tutz. 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, 2006.
- Jun Fan, Ting Hu, Qiang Wu, and Ding-Xuan Zhou. Consistency analysis of an empirical minimum error entropy algorithm. *Applied and Computational Harmonic Analysis*, 41(1):164–189, 2016.
- Yunlong Feng and Qiang Wu. Learning under $(1 + \epsilon)$ -moment conditions. *Submitted*, 2019.
- Yunlong Feng and Yiming Ying. Learning with correntropy-induced losses for regression with mixture of symmetric stable noise. *Applied and Computational Harmonic Analysis*, 48(2):795–810, 2019.
- Yunlong Feng, Xiaolin Huang, Lei Shi, Yuning Yang, and Johan A.K. Suykens. Learning with the maximum correntropy criterion induced losses for regression. *Journal of Machine Learning Research*, 16:993–1034, 2015.
- Frédéric Ferraty, Ali Laksaci, and Philippe Vieu. Functional time series prediction via conditional mode estimation. *Comptes Rendus Mathématique*, 340(5):389–392, 2005.
- Keinosuke Fukunaga and Larry Hostetler. The estimation of the gradient of a density function, with applications in pattern recognition. *IEEE Transactions on Information Theory*, 21(1):32–40, 1975.
- Ali Gannoun, Jerome Saracco, and Keming Yu. On semiparametric mode regression estimation. *Communications in Statistics - Theory and Methods*, 39(7):1141–1157, 2010.
- Ran He, Bao-Gang Hu, Wei-Shi Zheng, and Xiang-Wei Kong. Robust principal component analysis based on maximum correntropy criterion. *IEEE Transactions on Image Processing*, 20(6):1485–1494, 2011.
- Ran He, Tieniu Tan, Liang Wang, and Wei-Shi Zheng. $\ell_{2,1}$ -regularized correntropy for robust feature selection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2012*, pages 2504–2511. IEEE, 2012.
- Claudio Heinrich. The mode functional is not elicitable. *Biometrika*, 101(1):245–251, 2013.
- Eva Herrmann and Klaus Ziegler. Rates of consistency for nonparametric estimation of the mode in absence of smoothness assumptions. *Statistics & Probability Letters*, 68(4): 359–368, 2004.
- Ting Hu, Jun Fan, Qiang Wu, and Ding-Xuan Zhou. Learning theory approach to minimum error entropy criterion. *Journal of Machine Learning Research*, 14:377–397, 2013.
- Gordon C.R. Kemp and Joao Santos Silva. Regression towards the mode. *Journal of Econometrics*, 170(1):92–101, 2012.

- Myoung-Jae Lee. Mode regression. *Journal of Econometrics*, 42(3):337–349, 1989.
- Myoung-Jae Lee. Quadratic mode regression. *Journal of Econometrics*, 57(1):1–19, 1993.
- Myoung-Jae Lee and Hyun Ah Kim. Semiparametric econometric estimators for a truncated regression model: a review with an extension. *Statistica Neerlandica*, 52(2):200–225, 1998.
- Weifeng Liu, Puskal P. Pokharel, and Jose C. Principe. Correntropy: properties and applications in non-Gaussian signal processing. *IEEE Transactions on Signal Processing*, 55(11):5286–5298, 2007.
- Canyi Lu, Jinhui Tang, Min Lin, Liang Lin, Shuicheng Yan, and Zhouchen Lin. Correntropy induced l2 graph for robust subspace clustering. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1801–1808, 2013.
- Fusheng Lv and Jun Fan. Optimal learning with Gaussians and correntropy loss. *Analysis and Applications*, 2019. doi: 10.1142/S0219530519410124.
- Zhike Lv, Huiming Zhu, and Keming Yu. Robust variable selection for nonlinear models with diverging number of parameters. *Statistics & Probability Letters*, 91:90–97, 2014.
- Eric Matzner-Løfber, Ali Gannoun, and Jan G. De Gooijer. Nonparametric forecasting: a comparison of three kernel-based methods. *Communications in Statistics - Theory and Methods*, 27(7):1593–1617, 1998.
- Elias Ould-Saïd. A note on ergodic processes prediction via estimation of the conditional mode function. *Scandinavian Journal of Statistics*, 24(2):231–239, 1997.
- Emanuel Parzen. On estimation of a probability density function and mode. *The Annals of Mathematical Statistics*, 33(3):1065–1076, 1962.
- Karl F. Petty, Hisham Noeimi, Kumud Sanwal, Dan Rydzewski, Alexander Skabardonis, Pravin Varaiya, and Haitham Al-Deek. The freeway service patrol evaluation project: Database support programs, and accessibility. *Transportation Research Part C: Emerging Technologies*, 4(2):71–85, 1996.
- Jose C. Principe. *Information Theoretic Learning: Renyi’s Entropy and Kernel Perspectives*. Springer Science & Business Media, 2010.
- Alejandro Quintela-Del-Rio and Philippe Vieu. A nonparametric conditional mode estimate. *Journal of Nonparametric Statistics*, 8(3):253–266, 1997.
- Tim Robertson and Jonathan D. Cryer. An iterative procedure for estimating the mode. *Journal of the American Statistical Association*, 69(348):1012–1016, 1974.
- Thomas W. Sager and Ronald A. Thisted. Maximum likelihood estimation of isotonic modal regression. *The Annals of Statistics*, 10(3):690–707, 1982.
- Khaldani Salah and Yao Anne Françoise. Nonlinear parametric mode regression. *Communications in Statistics - Theory and Methods*, 46(6):3006–3024, 2016.

- Mrityunjay Samanta and Aerambamoorthy Thavaneswaran. Non-parametric estimation of the conditional mode. *Communications in Statistics-Theory and Methods*, 19(12):4515–4524, 1990.
- Hiroaki Sasaki, Yurina Ono, and Masashi Sugiyama. Modal regression via direct log-density derivative estimation. In *International Conference on Neural Information Processing*, pages 108–116. Springer, 2016.
- Ingo Steinwart and Andreas Christmann. *Support Vector Machines*. Springer, New York, 2008.
- Matt P. Wand and Chris M. Jones. *Kernel Smoothing*. Chapman & Hall, London, 1994.
- Xueqin Wang, Yunlu Jiang, Mian Huang, and Heping Zhang. Robust variable selection with exponential squared loss. *Journal of the American Statistical Association*, 108(502):632–643, 2013.
- Qiang Wu, Yiming Ying, and Ding-Xuan Zhou. Multi-kernel regularized classifiers. *Journal of Complexity*, 23(1):108–134, 2007.
- Weixin Yao and Longhai Li. A new regression model: modal linear regression. *Scandinavian Journal of Statistics*, 41(3):656–671, 2014.
- Weixin Yao and Sijia Xiang. Nonparametric and varying coefficient modal regression. *arXiv preprint arXiv:1602.06609*, 2016.
- Weixin Yao, Bruce G. Lindsay, and Runze Li. Local modal regression. *Journal of Nonparametric Statistics*, 24(3):647–663, 2012.
- Keming Yu and Katerina Aristodemou. Bayesian mode regression. *arXiv preprint arXiv:1208.0579*, 2012.
- Keming Yu, Katerina Aristodemou, Frauke Becker, and Joann Lord. Fast mode regression in big data analysis. In *Proceedings of the 2014 International Conference on Big Data Science and Computing*, page 24. ACM, 2014.
- Tong Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *The Annals of Statistics*, 32(1):56–85, 2004.
- Haiming Zhou and Xianzheng Huang. Nonparametric modal regression in the presence of measurement error. *Electronic Journal of Statistics*, 10(2):3579–3620, 2016.
- Haiming Zhou and Xianzheng Huang. Bandwidth selection for nonparametric modal regression. *Communications in Statistics-Simulation and Computation*, 48(4):968–984, 2019.