

Local Regression Smoothers with Set-Valued Outcome Data

Qiyu Li^a, Ilya Molchanov^b, Francesca Molinari^c, Sida Peng^d

^a*Department of Mathematical Statistics and Actuarial Science, University of Bern, and
Swiss Group for Clinical Cancer Research (SAKK)*

^b*Department of Mathematical Statistics and Actuarial Science, University of Bern*

^c*Department of Economics, Cornell University*

^d*Microsoft Research*

Abstract

This paper proposes a method to conduct local linear regression smoothing in the presence of set-valued outcome data. The proposed estimator is shown to be consistent, and its mean squared error and asymptotic distribution are derived. A method to build error tubes around the estimator is provided, and a small Monte Carlo exercise is conducted to confirm the good finite sample properties of the estimator. The usefulness of the method is illustrated on a novel dataset from a clinical trial to assess the effect of certain genes' expressions on different lung cancer treatments outcomes.

Keywords: Local regression smoothers; set valued outcome data; random sets; support function

1. Introduction

Statistical analysis has traditionally contended with problems of data imprecision due to limits in the measuring instruments and to measurement error, as well as with missing data, data coarsening and grouping. Geostatistical analysis
5 and mathematical morphology have contended with observational frameworks

*We are grateful to the Editor, the Area Editor, and two anonymous referees for comments that helped us substantially improve the paper. Molinari gratefully acknowledges support from NSF grant SES1824375.

*Corresponding author

Email address: fm72@cornell.edu (Francesca Molinari)

where the outcome of interest is a two or three dimensional set-valued object, e.g. a tumor or a grain. The common denominator of these challenging data-frameworks is the presence of set-valued data. Within the social sciences in particular, collection of data in the form of sets, especially intervals, has be-
10 come increasingly widespread. For example, the Health and Retirement Study is one of the first surveys where, in order to reduce item nonresponse, income data is collected from respondents in the form of brackets, with degenerate (singleton) intervals for individuals who opt to fully report their income (see, e.g. [1]). To reduce response burden, the Occupational Employment Statistics
15 (OES) program at the Bureau of Labor Statistics collects wage data from employers as intervals, and uses these data to construct estimates for wage and salary workers in 22 major occupational groups and 801 detailed occupations. Privacy concerns often motivate providing public use tax data as the number of tax payers in each of a finite number of cells. In the medical field, due to ethical
20 and cost reasons, time-to-event measurements are not collected on a continuous scale, but at pre-specified time intervals.

The partial identification literature in econometrics (e.g., [2]) has addressed the question of what can be learned about functionals of probability distributions, when some of the variables are only known to belong to (random) sets
25 and no assumptions are imposed on the distribution of the true variables within these sets. We take the identification results of this literature as our point of departure. Our contribution is to provide statistical results on local linear regression smoothing when the outcome data is set-valued and the regressors are exactly measured. Our paper relaxes the textbook setting (e.g., [3]) of nonpara-
30 metric regression – where regressors and outcome data $(\mathbf{x}_i, \mathbf{y}_i)$, $i = 1, \dots, n$, are precisely measured – by assuming that \mathbf{y}_i is only known to belong to an observed set \mathbf{Y}_i . In other words, we deal with an independently and identically distributed sample of observations for the pair $(\mathbf{x}_i, \mathbf{Y}_i)$ composed of a random vector \mathbf{x}_i in \mathbb{R}^m and a random convex compact set \mathbf{Y}_i in \mathbb{R}^d . Independence and
35 identical distribution for random sets and measurability of \mathbf{Y} are notions made precise in Appendix D, while in Section 2 we explain that the distribution of

\mathbf{Y} can be characterized as a *belief function*. The true (however unobservable) outcome associated with \mathbf{x} is a random vector \mathbf{y} that almost surely takes values in \mathbf{Y} . Our goal is to provide a nonparametric regression estimator for the
40 expectation conditional on \mathbf{x} of each random vector $\mathbf{y} \in \mathbf{Y}$. One can think of such expectation as the first-order moment of the belief function generated by \mathbf{Y} conditional on \mathbf{x} .

For a given tuple (\mathbf{x}, \mathbf{y}) that almost surely belongs to $\{\mathbf{x}\} \times \mathbf{Y}$, we denote by $m(x) = \mathbf{E}[\mathbf{y}|\mathbf{x} = x]$ the regression function for the chosen (\mathbf{x}, \mathbf{y}) . Each choice of $(\mathbf{x}, \mathbf{y}) \in \{\mathbf{x}\} \times \mathbf{Y}$ a.s. gives rise to a function m and we denote by \mathcal{M} the family of all regression functions generated in this manner. We let $M(x) \equiv \{m(x) : m \in \mathcal{M}\}$ and we observe that

$$M(x) = \mathbf{E}[\mathbf{Y}|\mathbf{x} = x] = \left\{ \mathbf{E}[\mathbf{y}|\mathbf{x} = x] : \mathbf{y} \in \mathbf{Y} \text{ a.s.} \right\}$$

is the conditional selection expectation of \mathbf{Y} , see [4, Sec. 2.1.6] and Section 2.

For example, consider the empirically relevant case that $d = 1$ and $\mathbf{Y} = [\mathbf{y}_L, \mathbf{y}_U]$ for two random variables $\mathbf{y}_L, \mathbf{y}_U$ such that $\mathbf{P}(\mathbf{y}_L \leq \mathbf{y}_U) = 1$. Then

$$M(x) = \left[\mathbf{E}[\mathbf{y}_L|\mathbf{x} = x], \mathbf{E}[\mathbf{y}_U|\mathbf{x} = x] \right]. \quad (1)$$

Our proposal is to estimate $M(x)$ as a weighted sum of the sets $\mathbf{Y}_1, \dots, \mathbf{Y}_n$,
45 with weights defined as in the local linear estimation literature.¹ The development of our technical results directly builds on classic references such as [5] and [6], and is closely related to [7] and [3].

For the case that $d = 1$, inspection of equation (1) might suggest to report an estimator given by the interval between a local constant or local linear regression
50 of \mathbf{y}_L on \mathbf{x} and one of \mathbf{y}_U on \mathbf{x} . Alternatively, it might suggest to report a local constant or local linear regression of the interval midpoint, $\tilde{\mathbf{y}} = (\mathbf{y}_L + \mathbf{y}_U)/2$, and of the interval width, $\mathbf{w} = \mathbf{y}_U - \mathbf{y}_L$, on \mathbf{x} . While both in finite sample

¹We comment on the case of local constant (Nadaraya–Watson) estimator in Appendix C.

and asymptotically these approaches are equivalent to what we propose for the case of a local constant regression, for the case of local linear regression equivalence breaks down in finite sample. The difference is important: we show in Remark 3.1 below that the alternative estimators just described may lead to a finite sample bias understating the width of $M(x)$ and are therefore unpalatable. For example, such estimators might be empty or a singleton in finite sample even though $M(x)$ is an interval of strictly positive width in population. In contrast, the estimator that we propose does not suffer from this problem, although it does have an asymptotic bias term similar to that of point identified local linear regression estimators.

Our approach is the first contribution in the literature to local regression smoothing when the set-valued outcome variable is in \mathbb{R}^d with $d > 1$. We derive the asymptotic properties of our estimator and extend results from [8] to obtain pointwise confidence bands that asymptotically cover the functional of interest with probability $1 - \alpha$. We report the results of Monte Carlo simulations with interval-valued \mathbf{Y} and with \mathbf{Y} being a ball randomly placed on the plane that support our theoretical findings.

We also demonstrate the usefulness of our approach with an empirical illustration that uses a novel dataset from a clinical trial on non-small-cell lung cancer patients, to study the relationship between tumor time to progression and specific gene expression measures.

Related literature. Within the partial identification literature, there is a large body of work analyzing regression with interval-valued data. [9] consider models where one variable (either outcome or covariate) is observed as intervals and all others are perfectly measured, and provide identification results for non-parametric as well as parametric models in this setting. [8] introduce to the partial identification literature the use of random set theory and provide results on identification and inference on best linear prediction parameters (ordinary least squares) when the outcome variable is interval-valued and the regressors are perfectly measured. [10] extend the familiar Sargan test for overidentifying

restrictions to the setting studied by [8]. [11] extend [8]’s approach to cover best linear approximation of any function $f(x)$ that is known to lie within two
85 identified bounding functions. [12] proposes an estimator for weighted average derivatives of conditional mean and conditional quantile functionals when either the outcome variable or a regressor is interval-valued. [13] propose empirical likelihood methods for random sets to conduct inference in the class of problems analyzed by [8]. All these papers focus exclusively on the case that
90 the set-valued outcome data is in \mathbb{R} .

In contrast, our approach leverages the theory of random sets to propose a set-valued local linear regression estimator for conditional set-valued expectations with $\mathbf{Y} \subset \mathbb{R}^d, d \geq 1$, and to establish its asymptotic properties. This estimator is novel in the literature, and so are our results establishing its consistency and asymptotic distribution.
95

The method that we propose differs significantly from other approaches in the statistical literature; see [14] for a discussion bridging this literature with partial identification. In particular, our proposal is distinct from the large and closely related literature that posits parametric models for set-valued data. In
100 these models tools from interval arithmetic are used to build analogs of the classic linear regression model for perfectly measured data, e.g. by assuming that $\mathbf{E}[\mathbf{Y}_i|\mathbf{x}_i] = A\mathbf{x}_i + B$, where A and B are intervals. See e.g. [15], [16], [17], and [18] among others for a discussion of least squares analysis of this and related models. [19] proposes nonparametric smoothing for this model, by
105 applying weighted least squares to the interval data and then using the resulting intercept as the estimator. [20] discuss various interpretations of set-valued data. Compared to this literature, we leave the conditional set-valued expectation completely unspecified, and nonparametrically estimate all regression functions compatible with the interval-valued data.

Finally, our proposal is distinct from the literature on data coarsening, e.g.
110 [21], [22] and [23]. In that literature, the key assumption of “coarsening at random” requires that for any possible value A of the random set \mathbf{Y} and a random vector \mathbf{y} that almost surely belongs to \mathbf{Y} , the conditional probability

$\mathbf{P}(\mathbf{Y} = A | \mathbf{y} = y_0)$ does not depend on $y_0 \in A$. This assumption restricts
115 directly the conditional distribution of the random set \mathbf{Y} , whereas we leave this
distribution completely unrestricted.

Structure of the paper. In Section 2 we set up our notation and we briefly review
local linear regression with singleton data. Our method implicitly applies it to
each tuple $(\mathbf{x}, \mathbf{y}) : (\mathbf{x}, \mathbf{y}) \in \{\mathbf{x}\} \times \mathbf{Y}$ a.s. In Section 3 we propose our estimator
120 and in Section 4 derive its asymptotic properties. In Section 5 we describe a
cross-validation method for bandwidth selection, and we extend the methods
proposed by [8] to test a hypothesis about the conditional expectation (eval-
uated at x_0) and to build pointwise error bands with prespecified asymptotic
coverage. In Section 6 we report the results of Monte Carlo experiments and
125 in Section 7 the results of our empirical illustration. Section 8 concludes. All
technical proofs are collected in Appendix A. Throughout we consider the case
that the regressors \mathbf{x} are random variables (random design case). In keeping
with the tradition in the statistics literature (e.g., [3]), we also report in Ap-
pendix B the case of deterministic design (nonstochastic explanatory variables).
130 Appendix C briefly discusses the local constant regression case. Appendix D
reports some basic facts in convex geometry and random set theory that we use
throughout the paper. We refer to [4] for a thorough account of random sets
theory. Appendix E provides additional simulation results.

2. Notation and preliminaries

135 We begin with listing our notation. We use boldface capital letters $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$
to denote random compact convex sets, normal font capital letters X, Y, Z and
 A, B, C to denote deterministic compact convex sets, boldface lower case letters
 $\mathbf{x}, \mathbf{y}, \mathbf{z}$ to denote random vectors or random variables, and normal font lowercase
letters x, y, z to denote deterministic vectors. For $x \in \mathbb{R}$, we denote the positive
140 and negative parts of x respectively by $x^+ = \max(0, x)$ and $x^- = -\min(0, x)$.
We let $(\Omega, \mathfrak{F}, \mathbf{P})$ denote a nonatomic probability space on which all random
vectors and random sets that we work with are defined, where Ω is the space of

elementary events equipped with σ -algebra \mathfrak{F} and probability measure \mathbf{P} . We denote the Euclidean space by \mathbb{R}^d , and equip it with the Euclidean norm (which is denoted by $\|\cdot\|$). We denote by $\mathcal{K}(\mathbb{R}^d)$ the collection of compact subsets of \mathbb{R}^d and by $\mathcal{K}_C(\mathbb{R}^d)$ the family of non-empty compact convex sets, also called convex bodies. We let $\mathbb{S}^{d-1} = \{x \in \mathbb{R}^d : \|x\| = 1\}$ denote the unit sphere in \mathbb{R}^d .

We assume that \mathbf{Y} is a random set in \mathbb{R}^d taking almost surely compact and convex values. In terms of measurability requirements, this amounts to

$$\{\omega : \mathbf{Y}(\omega) \cap K \neq \emptyset\} \in \mathfrak{F} \quad \forall K \in \mathcal{K}(\mathbb{R}^d). \quad (2)$$

The probabilities $\mathbf{P}(\mathbf{Y} \subseteq K)$, $K \in \mathcal{K}(\mathbb{R}^d)$, called the *containment functional* of \mathbf{Y} , fully characterize the distribution of \mathbf{Y} , [e.g., 4, Thm. 1.8.9]. As function of K , these probabilities are special cases of the *belief functions*, see [24] and more recently [25] and [26]. While general belief functions do not necessarily satisfy regularity conditions specific for the containment functional, the containment functionals are exactly semicontinuous belief functions. Then \mathbf{Y} describes the possible regions where a true value lies, and hence represents the ambiguity embedded in the observations, and coincides with the multivalued mapping Γ in [24].

To set the stage for local regression smoothing, we recall the standard construction of the local polynomial estimators for singleton-valued outcomes, see e.g. [6]. Suppose one is interested in estimating $\mathbf{E}(\mathbf{y}_i | \mathbf{x}_i = x_0)$ based on observations $(\mathbf{x}_i, \mathbf{y}_i)$, $i = 1, \dots, n$, where x_0 is a given value on the support of \mathbf{x} (e.g., a particular level of the gene expression measure in our empirical study). Then one fits a p -th order local model

$$\mathbf{y}_i = \theta_0(x_0) + \theta_1(x_0)(\mathbf{x}_i - x_0) + \dots + \theta_p(x_0)(\mathbf{x}_i - x_0)^p + \varepsilon_i,$$

using the regressor $\mathbf{x}_i - x_0$ (rather than \mathbf{x}_i) so that the intercept equals $\mathbf{E}(\mathbf{y}_i | \mathbf{x}_i = x_0)$. In this expression, the coefficients θ are written as a function of x_0 to emphasize that they change with the evaluation point (and this is what makes the

model “local”); to simplify notation, such dependence is suppressed henceforth. The local polynomial estimator of order p is then obtained by minimizing the weighted least squares

$$\sum_{i=1}^n \left(\mathbf{y}_i - \theta_0 - \theta_1(\mathbf{x}_i - x_0) - \cdots - \theta_p(\mathbf{x}_i - x_0)^p \right)^2 K\left(\frac{\mathbf{x}_i - x_0}{h_n}\right) \quad (3)$$

with respect to $\theta_0, \dots, \theta_p$. The kernel function $K(\cdot)$ is a nonnegative integrable function and the tuning parameter h_n is the bandwidth. As it is typically done, we assume that $h_n \rightarrow 0$ and $nh_n \rightarrow \infty$ as $n \rightarrow \infty$. The following condition on
160 the kernel function is imposed throughout this paper.

Assumption A (Kernel function). *The kernel $K(z)$, $z \in \mathbb{R}$, is a nonnegative function bounded above by $K_{\max} < \infty$, with compact support $[-c_K, c_K]$ for some $c_K \in (0, \infty)$, and satisfying*

$$\int K(z) dz = 1, \quad \int zK(z) dz = 0.$$

Denote $\text{Var}_K = \int z^2 K(z) dz$.

The normalization conditions on K are standard, while the compact support ensures that observations sufficiently far (compared to the order of the bandwidth) from the current point do not influence the estimator at this point, see
165 also Appendix B.

Solving explicitly the weighted least squares minimization problem in (3), one obtains the minimizer $\hat{\theta}$, and the first entry of it, the intercept $\hat{\theta}_0$, is used to estimate $m(x_0)$. This estimator can be written as

$$\hat{\mathbf{m}}(x_0) = \sum_{i=1}^n \ell_i(x_0) \mathbf{y}_i, \quad (4)$$

where

$$\begin{aligned}\ell_i(x_0) &= \frac{1}{nh_n} u^\top(0) \mathcal{B}_{nx_0}^{-1} u\left(\frac{\mathbf{x}_i - x_0}{h_n}\right) \boldsymbol{\kappa}_{in}, \\ u(z) &= \left(1, z, z^2/2!, \dots, z^p/p!\right)^\top, \\ \mathcal{B}_{nx_0} &= \frac{1}{nh_n} \sum_{i=1}^n u\left(\frac{\mathbf{x}_i - x_0}{h_n}\right) u^\top\left(\frac{\mathbf{x}_i - x_0}{h_n}\right) \boldsymbol{\kappa}_{in},\end{aligned}$$

with $\boldsymbol{\kappa}_{in} = K\left(\frac{\mathbf{x}_i - x_0}{h_n}\right)$. Note that $\ell_i(x_0)$, $i = 1, \dots, n$, sum up to one, and write

$$\mathbf{s}_j = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\kappa}_{in} (\mathbf{x}_i - x_0)^j, \quad j = 0, 1, \dots$$

It is easy to see that $\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 \geq 0$, and that the right-hand side of (4) is linear in the response variables, since the weights do not depend on the \mathbf{y}_i 's.

If $p = 0$ (local constant regression), $\hat{\mathbf{m}}(x_0)$ is the Nadaraya-Watson estimator with $\ell_i(x_0) = \boldsymbol{\kappa}_{in}/(n\mathbf{s}_0)$. If $p = 1$ (local linear regression), then

$$\ell_i(x_0) = \frac{\boldsymbol{\kappa}_{in}}{n} \frac{\mathbf{s}_2 - (\mathbf{x}_i - x_0) \mathbf{s}_1}{\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2}. \quad (5)$$

Our goal is to extend the local linear regression framework to set-valued outcomes: we propose an analog to estimator (4) with $p = 1$ and $\ell_i(x_0)$ as
170 given in (5), for the case that instead of knowing the exact value of \mathbf{y} , it is only assumed that \mathbf{y} almost surely belongs to a random set \mathbf{Y} . In this case \mathbf{y} is said to be a (measurable) *selection* of \mathbf{Y} . Distributions of all selections of \mathbf{Y} can be identified with the probability measures from the core of the belief function generated by \mathbf{Y} , that is, probability measures dominating the belief function.
175 The pair (\mathbf{x}, \mathbf{y}) is a selection of $\{\mathbf{x}\} \times \mathbf{Y}$, a random closed set in $I \times \mathbb{R}^d$ with I the support of \mathbf{x} . This framework can alternatively be described as associating with each value of the explanatory variable \mathbf{x} a belief function describing the (conditional) distribution of \mathbf{Y} .

Whereas in the standard case of singleton-valued outcomes one observes
180 singleton-valued data $(\mathbf{x}_i, \mathbf{y}_i)$, $i = 1, \dots, n$, in our framework the observations

are set-valued, $(\mathbf{x}_i, \mathbf{Y}_i)$, $i = 1, \dots, n$. As a result, our estimators are also set-valued, and in order to assess their properties, we need to define square loss for sets, so as to formalize consistency results and the notion of mean squared error. To do so, and to provide a computationally tractable estimator, we
185 exploit the duality between convex sets and their *support function* (see, e.g., Chapter 13 in [27], and (D.2) in Appendix D). The support function of \mathbf{Y} in direction $v \in \mathbb{S}^{d-1}$ is given by $s(\mathbf{Y}, v) \equiv \sup_{y \in \mathbf{Y}} v^\top y$, and can be used to define the *width function* of \mathbf{Y} in direction $v \in \mathbb{S}^{d-1}$, $w(\mathbf{Y}, v) \equiv s(\mathbf{Y}, v) + s(\mathbf{Y}, -v)$ (see Appendix D). We assume that \mathbf{Y} is integrably bounded, that is, $\|\mathbf{Y}\| =$
190 $\sup_{y \in \mathbf{Y}} \|y\|$ is integrable (Assumption B in the next section provides sufficient conditions guaranteeing that this is the case), and since $|s(\mathbf{Y}, v)| \leq \|\mathbf{Y}\|$ for all v from the unit sphere, this implies that the support function is integrable. It is possible to show that $\mathbf{E}s(\mathbf{Y}, v) = s(\mathbf{E}\mathbf{Y}, v)$ [see 4, Theorem 2.1.35], i.e. the expected support function is the support function of a convex body $\mathbf{E}\mathbf{Y}$, which
195 in turn is called the *expectation* of \mathbf{Y} . This expectation equals the set of values $\mathbf{E}\mathbf{y}$ for all random vectors \mathbf{y} such that $\mathbf{y} \in \mathbf{Y}$ a.s.

Similarly, for given x it is possible to define the *conditional expectation*

$$\mathbf{E}[\mathbf{Y}|\mathbf{x} = x] = \left\{ \mathbf{E}[\mathbf{y}|\mathbf{x} = x] : \mathbf{y} \in \mathbf{Y} \text{ a.s.} \right\},$$

and also in this case it holds that $\mathbf{E}[s(\mathbf{Y}, v)|\mathbf{x} = x] = s(\mathbf{E}[\mathbf{Y}|\mathbf{x} = x], v)$ [see, e.g., 4, Sec. 2.1.6]. The set $\mathbf{E}[\mathbf{Y}|\mathbf{x} = x]$ is the object of interest in this paper, and one can think of it as the first-order moment of the belief function generated by
200 \mathbf{Y} conditional on \mathbf{x} .

To simplify the exposition, henceforth we assume that \mathbf{x} is a scalar random variable and that I is an interval, $I \subset \mathbb{R}$. Our results apply, subject only to modification in notation and convergence rates (as in the point identified case), with vector-valued \mathbf{x} provided the real-valued bandwidth is replaced by
205 a matrix-valued one.

The family of support functions of all non-empty compact convex subsets in \mathbb{R}^d is a subset of the family of continuous functions on the unit sphere \mathbb{S}^{d-1} . In

particular, the Hausdorff metric between compact convex sets equals the uniform (L_∞) distance between their support functions, see e.g. [28, Lemma 1.8.14]. For our purposes, it is convenient to endow the family of continuous functions on the unit sphere with the L_2 -metric, so that the distance between two non-empty compact convex sets A_1 and A_2 is given by

$$L(A_1, A_2) = \left(\int_{\mathbb{S}^{d-1}} (s(A_1, v) - s(A_2, v))^2 dv \right)^{\frac{1}{2}}. \quad (6)$$

The integration is performed with respect to the uniform measure on \mathbb{S}^{d-1} . If $d = 1$, the integral turns into the sum of two terms for $v = 1$ and $v = -1$. The distance to the empty set is assigned to be infinite.

In Section 3, we employ this distance to define the mean square error of our estimator. This distance differs from the standard Hausdorff distance used in the related literature in partial identification and in the standard laws of large numbers and central limit theorems for Minkowski averages of random sets. However, under our assumptions the result of Theorem 3 in [29] yields that these two metrics define the same topology, and so the consistency with respect to the L_2 -distance implies consistency with respect to the L_∞ -distance. At the same time, use of the L_2 -distance is particularly well suited to analyze properties of estimators based on least squares minimization.

3. Nonparametric smoothing for random sets

In the following we assume that $(\mathbf{x}_i, \mathbf{Y}_i)$, $i = 1, \dots, n$, is a sample of i.i.d. realizations of (\mathbf{x}, \mathbf{Y}) as defined in Appendix D, where \mathbf{Y} satisfies Assumption B introduced below. This i.i.d. assumption is consistent with many collection processes of set-valued data, such as, e.g., the use of unfolding brackets in the Health and Retirement Study, in the Occupational Employment Statistics survey of the Bureau of Labor Statistics, and in the empirical application that we present in Section 7. We relate it to the typical i.i.d. assumption for singleton-valued data following our statement of Assumption B below.

When the outcome data is set-valued, it is necessary to obtain an estimator for the collection of conditional expectations $\mathbf{E}[\mathbf{y}|\mathbf{x} = x]$ for all $(\mathbf{x}, \mathbf{y}) \in \{\mathbf{x}\} \times \mathbf{Y}$ a.s. This can be accomplished by repeating the procedure in the previous section for all selections of $\{\mathbf{x}\} \times \mathbf{Y}$. Computationally this is easily achieved by taking the weighted Minkowski average of the \mathbf{Y}_i data (see Appendix D for a formal definition of Minkowski sum):

$$\hat{\mathbf{M}}(x_0) = \sum_{i=1}^n \ell_i(x_0) \mathbf{Y}_i. \quad (7)$$

For $p = 0$ we obtain a local constant set-valued regression estimator; the choice $p = 1$ yields a local linear set-valued regression estimator. Note that (7) is also the Fréchet mean of the observed values $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ in the metric given by (6),
230 see [30] and Sec. 2.2.5 in [4].

The estimator in (7) yields a convex set, therefore we can characterize its properties by working with its support function (see (D.2) in Appendix D and Chapter 13 of [27]). To simplify notation, in what follows we omit the argument x_0 in $\ell_i(x_0)$ and write shortly ℓ_i , unless the dependence on x_0 is essential. By representing the difference of its positive and negative parts as $\ell_i = \ell_i^+ - \ell_i^-$, and using that $s(-A, v) = s(A, -v)$ for a convex compact set A and its centrally symmetric set $-A = \{-x : x \in A\}$, we arrive at

$$\begin{aligned} s(\hat{\mathbf{M}}(x_0), v) &= s\left(\sum_{i=1}^n (\ell_i^+ - \ell_i^-) \mathbf{Y}_i, v\right) = \sum_{i=1}^n \ell_i^+ s(\mathbf{Y}_i, v) + \sum_{i=1}^n \ell_i^- s(\mathbf{Y}_i, -v) \\ &= \sum_{i=1}^n (\ell_i + \ell_i^-) s(\mathbf{Y}_i, v) + \sum_{i=1}^n \ell_i^- s(\mathbf{Y}_i, -v) = \sum_{i=1}^n \ell_i s(\mathbf{Y}_i, v) + \sum_{i=1}^n \ell_i^- w(\mathbf{Y}_i, v). \end{aligned}$$

A key feature of the above estimator is that it averages the support function of the set \mathbf{Y}_i in direction $+v$ when $\ell_i > 0$, and in direction $-v$ when $\ell_i < 0$. In doing so we guarantee that the estimator is always *non-empty* for any n , a highly desirable feature in light of Assumption B.

Remark 3.1. When $d = 1$ and $\mathbf{Y} = [\mathbf{y}_L, \mathbf{y}_U]$ with $\mathbf{P}(\mathbf{y}_U \geq \mathbf{y}_L) = 1$, one might

consider two estimators as alternatives to $\hat{M}(x_0)$. One is given by

$$\hat{N}(x_0) = \left[\sum_{i=1}^n \ell_i \mathbf{y}_{iL}, \sum_{i=1}^n \ell_i \mathbf{y}_{iU} \right].$$

The other is obtained by regressing the midpoint ($\tilde{\mathbf{y}}$) and the width (\mathbf{w}) of the interval $[\mathbf{y}_L, \mathbf{y}_U]$ on \mathbf{x} and letting

$$\hat{O}(x_0) = \left[\sum_{i=1}^n \ell_i \tilde{\mathbf{y}}_i - \sum_{i=1}^n \ell_i \frac{\mathbf{w}_i}{2}, \sum_{i=1}^n \ell_i \tilde{\mathbf{y}}_i + \sum_{i=1}^n \ell_i \frac{\mathbf{w}_i}{2} \right].$$

Standard arguments in [5] yield that $\hat{N}(x_0)$ and $\hat{O}(x_0)$ are consistent estimators of

$$M(x_0) = \mathbf{E}[\mathbf{Y}|\mathbf{x} = x_0] = \left[\mathbf{E}[\mathbf{y}_L|\mathbf{x} = x_0], \mathbf{E}[\mathbf{y}_U|\mathbf{x} = x_0] \right]$$

with respect to the L_2 -distance. However, these estimators can have large finite sample bias, and even be empty (with asymptotically vanishing probability), as illustrated in the following example. Suppose that for i with $\ell_i > 0$, $\mathbf{y}_{iL} = \mathbf{y}_{iU}$; and for i with $\ell_i < 0$, $\mathbf{y}_{iU} > \mathbf{y}_{iL}$.² Then

$$\begin{aligned} \sum_{i=1}^n \ell_i \mathbf{y}_{iL} &= \sum_{i=1}^n \ell_i^+ \mathbf{y}_{iL} - \sum_{i=1}^n \ell_i^- \mathbf{y}_{iL} = \sum_{i=1}^n \ell_i^+ \mathbf{y}_{iU} - \sum_{i=1}^n \ell_i^- \mathbf{y}_{iL} \\ &> \sum_{i=1}^n \ell_i^+ \mathbf{y}_{iU} - \sum_{i=1}^n \ell_i^- \mathbf{y}_{iU} = \sum_{i=1}^n \ell_i \mathbf{y}_{iU}, \end{aligned}$$

235 and $\hat{N}(x_0)$ is empty. One can similarly show that $\hat{O}(x_0)$ is empty. Similarly empty estimators may result even if $\mathbf{y}_{iU} > \mathbf{y}_{iL}$ whenever $\ell_i > 0$, depending on the realizations of \mathbf{y}_{iL} and \mathbf{y}_{iU} , see Figure 1 for $\hat{N}(x_0)$. Even if one censors $\mathbf{w}_i = 0$ if $\ell_i < 0$, the resulting estimator may still in finite sample significantly understate the width of $M(x_0)$.

240 While the example in Remark 3.1 might appear stylized, it highlights a

²While the example is provided for the case $d = 1$, similar constructions can be obtained also when $d \geq 2$.

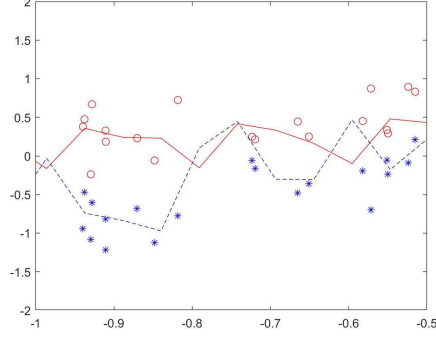


Figure 1: Possible emptiness of the estimator $\hat{N}(x_0)$. Blue dashed line: $\sum_{i=1}^n \ell_i \mathbf{y}_{iL}$; red solid line: $\sum_{i=1}^n \ell_i \mathbf{y}_{iU}$.

finite sample problem that can easily occur in practice with interval-valued data, but does not affect the corresponding estimators in the singleton-valued case. The reason is that in the singleton case, local regression smoothers are weighted averages of the observed outcomes. That is also the case for our
245 estimator, $\hat{M}(x_0)$, which averages the *sets* \mathbf{Y}_i and indeed is always non-empty. On the other hand, $\hat{N}(x_0)$ and $\hat{O}(x_0)$ average specific *selections* of \mathbf{Y}_i (e.g., the extreme points), without recognizing that the sign of the weight may affect which selection is extreme in a given direction.

Throughout the paper we assume $I = \mathbb{R}$ and we impose the following re-
250 strictions on the observed and theoretical responses and on the density function of \mathbf{x} .

- Assumption B** (Observed responses). (i) Let $(\mathbf{x}_i, \mathbf{Y}_i)$, $i = 1, \dots, n$, be a sample of i.i.d. realizations of (\mathbf{x}, \mathbf{Y}) , $i = 1, \dots, n$. Conditional on $\mathbf{x}_1, \dots, \mathbf{x}_n$, the observations $\mathbf{Y}_1, \dots, \mathbf{Y}_n$, are non-empty random compact convex sets.
255 (ii) $\mathbf{Y}_i \subset \xi_i + B$ a.s. for square integrable random vectors ξ_i , $i = 1, \dots, n$, and a deterministic compact set B that is the same for all i .

Define

$$\varepsilon_i(v) \equiv s(\mathbf{Y}_i, v) - s(M(\mathbf{x}_i), v), \quad v \in \mathbb{S}^{d-1}. \quad (8)$$

By Assumption B, $\varepsilon_i(\cdot)$, $i = 1, \dots, n$, are i.i.d. copies of a square integrable random function $\varepsilon(v)$, $v \in \mathbb{S}^{d-1}$, such that $\mathbf{E}[\varepsilon_i(v)|\mathbf{x}_i] = 0$ \mathbf{x}_i -a.s. for all $v \in \mathbb{S}^{d-1}$. The square integrability follows from the inequality,

$$\varepsilon_i(v) \leq s(B, v) + |\xi_i^\top v| + |\eta_i^\top v|,$$

where η_i is a square integrable selection of $M(\mathbf{x}_i)$. This selection exists in view of Assumption B(ii) and can be chosen as the point of $M(\mathbf{x}_i) = \mathbf{E}(\mathbf{Y}_i|\mathbf{x}_i) \subset \mathbf{E}(\xi_i|\mathbf{x}_i) + B$ nearest to $\mathbf{E}(\xi_i|\mathbf{x}_i)$. Note that ε does not admit a geometric interpretation as the support function of a random set.

Denote by $C(v, u) = \mathbf{E}[\varepsilon(v)\varepsilon(u)]$ the covariance function of ε and let σ_{\max}^2 be the supremum of $C(v, v) = \mathbf{E}[\varepsilon(v)^2]$ over all v from the unit sphere. Assumption B(ii) guarantees that \mathbf{Y}_i is uniformly integrably bounded, and implies that the diameters of all \mathbf{Y}_i 's are bounded by a deterministic constant. Hence, the ambiguity range is limited to belong to a deterministic set, and σ_{\max}^2 is finite.

It is worth to compare our random sampling assumption with the standard one for singleton-valued variables. In that context, one has $\mathbf{y}_i = m(\mathbf{x}_i) + \varepsilon_i$, and $(\mathbf{x}_i, \mathbf{y}_i)$ are assumed i.i.d., and as a consequence ε_i are i.i.d. In our context, we assume that $(\mathbf{x}_i, \mathbf{Y}_i)$ are i.i.d., and as a consequence $\varepsilon_i(v)$ are i.i.d.

In dimension $d = 1$, we have $s(\mathbf{Y}_i, 1) = \mathbf{y}_{iU}$, $s(\mathbf{Y}_i, -1) = -\mathbf{y}_{iL}$, and Part (i) of Assumption B requires that $\mathbf{y}_{iL} = \mathbf{E}[\mathbf{y}_L|\mathbf{x}] - \varepsilon_i(-1)$, $\mathbf{y}_{iU} = \mathbf{E}[\mathbf{y}_U|\mathbf{x}] + \varepsilon_i(1)$ with $\varepsilon_i(1) + \varepsilon_i(-1) \geq -(\mathbf{E}[\mathbf{y}_U|\mathbf{x}] - \mathbf{E}[\mathbf{y}_L|\mathbf{x}])$ almost surely. The latter condition replicates the requirement that $\mathbf{P}(\mathbf{y}_U \geq \mathbf{y}_L) = 1$.

Next, we require the conditional expectation of $\mathbf{E}[\mathbf{Y}|\mathbf{x}]$ to have a sufficiently smooth support function, thereby allowing for standard expansions used in obtaining the asymptotic properties of the local linear estimator.

Assumption C (Theoretical response function). *The function $M(x)$, $x \in \mathbb{R}$, is such that $s(M(x), v)$ admits a second derivative $s''(M(x), v)$ in x , uniformly bounded for all $v \in \mathbb{S}^{d-1}$.*

In dimension $d = 1$, Assumption C means the second order differentiability

of the end-points of the interval-valued function $M(x)$. Finally, we assume that the common density f of the independent design points satisfies the following condition, which is similar to that imposed in Condition 1(ii) of [5] with singleton responses. This is a standard condition in nonparametric regression; it
285 guarantees that the design points are not too concentrated in some areas.

Assumption D (Density). *The density f is strictly positive at x_0 and belongs to the family $\mathcal{H}(1, \gamma)$ of Lipschitz functions with constant $\gamma > 0$, that is,*

$$|f(x') - f(x'')| \leq \gamma |x' - x''|$$

for all $x', x'' \in \mathbb{R}$.

We measure the quality of $\hat{M}(x_0)$ as set-valued estimator of $M(x_0)$ by the quadratic loss function defined in (6),

$$L(\hat{M}(x_0), M(x_0))^2 = \int_{\mathbb{S}^{d-1}} (s(\hat{M}(x_0), v) - s(M(x_0), v))^2 dv.$$

The mean squared error (MSE) of the estimator is then the expectation of $L(\hat{M}(x_0), M(x_0))^2$. A classic bias-variance decomposition yields

$$\text{MSE}(x_0) = \int_{\mathbb{S}^{d-1}} b_{x_0}^2(v) dv + \int_{\mathbb{S}^{d-1}} \sigma_{x_0}^2(v) dv,$$

where $b_{x_0}^2(v)$ and $\sigma_{x_0}^2(v)$ are squared bias and variance, given by

$$\begin{aligned} b_{x_0}^2(v) &= \mathbf{E} \left(\mathbf{E}[s(\hat{M}(x_0), v) | \mathbf{x}_1, \dots, \mathbf{x}_n] - s(M(x_0), v) \right)^2, \\ \sigma_{x_0}^2(v) &= \mathbf{E} \left(s(\hat{M}(x_0), v) - s(\mathbf{E}[\hat{M}(x_0) | \mathbf{x}_1, \dots, \mathbf{x}_n], v) \right)^2. \end{aligned}$$

Because $\mathbf{E}[\mathbf{Y}_i | \mathbf{x}_i] = M(\mathbf{x}_i)$, we have

$$\mathbf{E}[s(\hat{M}(x_0), v) | \mathbf{x}_1, \dots, \mathbf{x}_n] = \sum_{i=1}^n \ell_i s(M(\mathbf{x}_i), v) + \sum_{i=1}^n \ell_i^- w(M(\mathbf{x}_i), v).$$

Rearranging the terms, we arrive at

$$b_{x_0}^2(v) = \mathbf{E} \left(\sum_{i=1}^n \ell_i (s(M(\mathbf{x}_i), v) - s(M(x_0), v)) + \sum_{i=1}^n \ell_i^- w(M(\mathbf{x}_i), v) \right)^2 \quad (9)$$

and

$$\sigma_{x_0}^2(v) = \mathbf{E} \left(\sum_{i=1}^n \ell_i (s(\mathbf{Y}_i, v) - s(M(\mathbf{x}_i), v)) + \sum_{i=1}^n \ell_i^- (w(\mathbf{Y}_i, v) - w(M(\mathbf{x}_i), v)) \right)^2.$$

By Assumption B, the variance can be expressed as

$$\sigma_{x_0}^2(v) = \mathbf{E} \left(\sum_{i=1}^n \ell_i \varepsilon_i(v) + \sum_{i=1}^n \ell_i^- (\varepsilon_i(v) + \varepsilon_i(-v)) \right)^2. \quad (10)$$

Differently from the classical case with singleton responses \mathbf{y}_i , the *negative* parts of the weights in (9) play an essential role with set-valued responses. This is because while the difference between $s(M(\mathbf{x}_i), v)$ and $s(M(x_0), v)$ is small when \mathbf{x}_i is close to x_0 , the width $w(M(\mathbf{x}_i), v)$ does not vanish as \mathbf{x}_i becomes closer to x_0 . Thus, the bias increases by a constant and may not tend to zero if some weights are negative and not close to zero. Much of our asymptotic analysis is concerned with establishing the asymptotic behavior of these negative weights.

The methodology that we propose for local linear regression smoothing can be applied also in the case of local polynomial regression models with $p \geq 2$. In this case, however, extra care is required to show that the negative weights are asymptotically negligible.

4. Asymptotic properties of the set-valued estimators

In the local linear regression setting, negative weights may appear in (9) and hence affect the bias in the case of set-valued data. Following [5], in order to avoid zero in the denominator of the local linear estimator, we redefine ℓ_i by letting

$$\ell_i = \frac{\kappa_{in}}{n} \frac{\mathbf{s}_2 - (\mathbf{x}_i - x_0)\mathbf{s}_1}{\mathbf{s}_2\mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4}}. \quad (11)$$

We use \mathcal{O} and \mathcal{O} to denote the deterministic order of magnitude uniformly in $f \in \mathcal{H}(1, \gamma)$. For a sequence $\{z_n, n \geq 1\}$ of random variables determined through the design points and the observations, write $z_n = \mathcal{O}_r(a_n)$ if

$$\sup_{f \in \mathcal{H}(1, \gamma)} \mathbf{E}|z_n|^r = \mathcal{O}(a_n^r).$$

The notation $\mathcal{O}_r(a_n)$ is defined similarly. We then have $\mathcal{O}_r(a_n)\mathcal{O}_r(b_n) = \mathcal{O}_{r/2}(a_nb_n)$, and

$$z_n = \mathbf{E}z_n + \mathcal{O}_r(\mathbf{E}|z_n - \mathbf{E}z_n|^r)^{1/r}.$$

To determine the contribution to the bias resulting from the negative weights, we first derive the expected sum of the squared weights ℓ_i^2 . Proofs of the following results are given in Appendix A.

Proposition 4.1. *Let $h_n \rightarrow 0$ and $nh_n \rightarrow \infty$ as $n \rightarrow \infty$. Under Assumptions A and D,*

$$\mathbf{E} \sum_{i=1}^n \ell_i^2 = \frac{1}{nh_n f(x_0)} \int K^2(z) dz + \mathcal{O}\left(\frac{1}{nh_n}\right). \quad (12)$$

Next, we obtain the second moment of the sum of the negative weights.

Proposition 4.2. *Let $h_n \rightarrow 0$ and $nh_n \rightarrow \infty$ as $n \rightarrow \infty$. Under Assumptions A and D, for sufficiently large r ,*

$$\mathbf{E} \left(\sum_{i=1}^n \ell_i^- \right)^2 = \frac{1}{h_n} \mathcal{O}\left((1/\sqrt{nh_n})^r\right).$$

With this result in hand, we can derive the mean squared error of our estimator. As the mean squared error converges to zero as n increases to infinity, this result yields consistency of our estimator as well as its rate of convergence.

Theorem 4.3. *Under Assumptions A, B, C, and D, if $h_n = cn^{-\beta}$ with $0 < \beta < 1$ and a constant $c > 0$, the mean squared error of the local linear estimator (7) is*

$$\text{MSE}(x_0) = \frac{h_n^4 (\text{Var}_K)^2}{4} \int_{\mathbb{S}^{d-1}} s''(M(x_0), v)^2 dv + \frac{\int_{\mathbb{S}^{d-1}} C(v, v) dv}{nh_n f(x_0)} \int K^2(z) dz + \mathcal{O}\left(h_n^4 + \frac{1}{nh_n}\right).$$

We conclude this section by deriving a limit theorem for the support function of the estimators as processes on the unit sphere. In turn, this limit theorem can be used to build error tubes for the estimator as explained in Section 5. Let $\zeta(v)$, $v \in \mathbb{S}^{d-1}$, be a centered Gaussian process on the unit sphere with the covariance

$$\mathbf{E}[\zeta(v)\zeta(u)] = \frac{C(v, u)}{f(x_0)} \int K(z)^2 dz. \quad (13)$$

Theorem 4.4. *Assume that $h_n = cn^{-\beta}$ with $0 < \beta < 1$, and fix $x_0 \in I$. Under Assumptions A, B, C, and D, the stochastic process*

$$\sqrt{nh_n} \left(s(\hat{M}(x_0), v) - s(M(x_0), v) - h_n^2 \frac{1}{2} s''(M(x_0), v) \text{Var}_K \right)$$

constructed using the local linear estimator in (7) converges in distribution in the space of continuous functions on \mathbb{S}^{d-1} with the uniform metric to the Gaussian process ζ .

5. Cross-validation and error tubes

Cross-validation. In the classical setting, where the observation pairs $(\mathbf{x}_i, \mathbf{y}_i)$ are real-valued, one typically chooses the bandwidth h_n to minimize the leave-one-out cross-validation score, defined as

$$\text{CV} = \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \hat{\mathbf{m}}_{(-i)}(\mathbf{x}_i))^2,$$

where $\hat{\mathbf{m}}_{(-i)}(x) = \sum_{j=1}^n \mathbf{y}_j \ell_{j,(-i)}(x)$ and

$$\ell_{j,(-i)}(x) = \begin{cases} 0 & \text{if } j = i, \\ \frac{\ell_j(x)}{\sum_{k \neq i} \ell_k(x)} & \text{if } j \neq i. \end{cases}$$

310 This procedure assigns weight zero to \mathbf{x}_i and renormalizes the other weights to sum to one.

Following the same idea, we define the cross-validation score for the set-

valued responses \mathbf{Y}_i as

$$\text{CV} = \frac{1}{n} \sum_{i=1}^n \int_{\mathbb{S}^{d-1}} (s(\mathbf{Y}_i, v) - s(\hat{\mathbf{M}}_{(-i)}(\mathbf{x}_i), v))^2 dv, \quad (14)$$

where $\hat{\mathbf{M}}_{(-i)}(x) = \sum_{j=1}^n \mathbf{Y}_j \ell_{j,(-i)}(x)$. If one is interested in a specific projection in direction v , the above expression simplifies by removing the integral.

If $\mathbf{Y}_i = [\mathbf{y}_{iL}, \mathbf{y}_{iU}] \subset \mathbb{R}$, (14) turns into

$$\text{CV} = \frac{1}{n} \sum_{i=1}^n \left(\mathbf{y}_{iL} - \hat{\mathbf{M}}_{(-iL)}(\mathbf{x}_i) \right)^2 + \left(\mathbf{y}_{iU} - \hat{\mathbf{M}}_{(-iU)}(\mathbf{x}_i) \right)^2, \quad (15)$$

where $\hat{\mathbf{M}}_{(-iL)}(\mathbf{x}_i)$ and $\hat{\mathbf{M}}_{(-iU)}(\mathbf{x}_i)$ denote the lower and upper bounds of $\hat{\mathbf{M}}_{(-i)}(\mathbf{x}_i)$.

315 We denote by $h_{n,\text{CV}}$ the bandwidth that minimizes (15) (or (14), depending on the application).

Error tubes. The optimal bandwidth which minimizes the MSE in Theorem 4.3 is $h_{n,\text{mse}} = Cn^{-1/5}$, with some constant C that does not depend on n . However, such a choice of bandwidth implies $nh_n^5 \not\rightarrow 0$ and the leading bias
320 term in Theorem 4.4 does not vanish, as in the classical case for singleton-valued outcomes. Similarly to that case, one can use undersmoothing as an approach to bias reduction. In Section 6 we illustrate the impact of undersmoothing on the error tubes that we describe next.

Rather than undersmooth, we propose to report statistical uncertainty in our estimates in the form of pointwise error tubes – an analog of error bands for singleton-valued data. Specifically, for each value x_0 considered we propose to report the set

$$\hat{\mathcal{C}}(x_0) = \hat{\mathbf{M}}(x_0) + \frac{c_\alpha}{\sqrt{nh_n}} B, \quad (16)$$

where $B = \{b : \|b\| \leq 1\}$ is the unit ball. In (16) c_α is chosen so that

$$\mathbf{P} \left(\max_{v: \|v\|=1} \{\zeta(v)\}_+ > c_\alpha \right) = \alpha, \quad (17)$$

where ζ is the centered Gaussian process with covariance kernel (13), see Theo-

rem 4.4. The critical value c_α can be obtained by simulation, or can be estimated using the bootstrap. Validity of the bootstrap can be formally established as in Proposition 2.1 of [8] [see also 31, Theorem 4.13]. It follows from Theorem 4.4 that

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(\max_{v: \|v\|=1} \{s(\hat{\mathbf{M}}(x_0), v) - s(M(x_0), v) - h_n^2 \frac{1}{2} s''(M(x_0), v) \text{Var}_K - s(\hat{\mathcal{C}}(x_0), v)\}_+ = 0 \right) \geq 1 - \alpha. \quad (18)$$

If one is interested in a specific projection in direction v , a valid error band for $s(M(x_0), v)$ is obtained by replacing (16) with

$$\left[s(\hat{\mathbf{M}}(x_0), v) - \frac{c_{\alpha, v}}{\sqrt{nh_n}}, s(\hat{\mathbf{M}}(x_0), v) + \frac{c_{\alpha, v}}{\sqrt{nh_n}} \right]. \quad (19)$$

where $c_{\alpha, v}$ is obtained as in (17) replacing the maximization over v with $\|v\| = 1$ by a fixed direction v .

Existing methods of bias correction (other than undersmoothing, the effect of which we are already investigating in our Monte Carlo exercise) could be extended to the case of set-valued outcomes. However, we do not report such findings here,³ because any form of bias reduction may result in an empty estimator, which we regard as an undesirable feature as discussed in Remark 3.1.

6. Monte Carlo Simulations

We perform a simulation study for the case that $d = 1$ and for the case that $d = 2$. In the first case, we use the following data generating process (DGP1):

$$\begin{aligned} \mathbf{y}_L &= 0.90 + 1.27\mathbf{x} + 5.18\mathbf{x}^2 - \varepsilon_L \\ \mathbf{y}_U &= 0.90 + 1.27\mathbf{x} + 10.18\mathbf{x}^2 + \varepsilon_U, \end{aligned}$$

³Although these are available from the authors upon request.

Table 1: Coverage probability at 95% nominal level using cross-validation for DGP1.

sample size	x_0	Coverage of $M(x_0)$	Coverage of $\mathbf{E}(\hat{M}(x_0))$	Coverage of $M(x_0)$ with $h = 1/2h_{n,CV}$	Coverage of $M(x_0)$ with $h = 1/3h_{n,CV}$
200	-0.4	0.8315	0.8245	0.9055	0.9695
	0	0.8855	0.8550	0.8565	0.9515
	0.2	0.9330	0.9270	0.9865	0.9980
	0.4	0.9270	0.9040	0.9255	0.9875
500	-0.4	0.8580	0.8485	0.9300	0.9790
	0	0.9245	0.9095	0.9710	0.9920
	0.2	0.9240	0.9200	0.9710	0.9950
	0.4	0.9340	0.9145	0.9180	0.9760
1000	-0.4	0.8910	0.8760	0.9430	0.9845
	0	0.9035	0.8935	0.9360	0.9830
	0.2	0.9230	0.9210	0.9570	0.9890
	0.4	0.9225	0.9125	0.9125	0.9760
2000	-0.4	0.88200	0.8710	0.9450	0.9835
	0	0.9020	0.8915	0.9390	0.9870
	0.2	0.9320	0.9125	0.9525	0.9900
	0.4	0.9335	0.9170	0.9635	0.9915

with \mathbf{x} drawn from a Beta distribution with support shifted to be $[-1, 1]$ and with shape parameters $(2, 4)$, and ε_L and ε_U drawn independently from a Uniform distribution on $[0, 1]$. We let the sample size $n = 200, 500, 1000, 2000$. For values of $x_0 = 0, 0.2, 0.4, 0.6$ we evaluate the coverage probability of the error tubes in equation (16).

We compare different implementations of the error tubes, and in Table 1 we report: (i) coverage probability of the true set $\mathbf{M}(x_0)$ by the error tube (meaning that the true set is a subset of the tube) in (16) computed using the cross-validation bandwidth (column 3); (ii) coverage probability as in (18), with the error tube in (16) computed using the cross-validation bandwidth (column 4); (iii) same exercise as in (i) but using undersmoothed bandwidths (columns 5 and 6). The results are based on 200 Monte Carlo replications.

In these simulations, the asymptotic bias does not affect the ability of the error tube in (16) to cover the true set $M(x_0)$ compared to $\mathbf{E}[\hat{\mathbf{M}}(x_0)]$, see columns (3) and (4) of the table. If we undersmooth the bandwidth, the confidence interval enlarges substantially and coverage of the true set becomes conservative.

In Appendix E we report the results of two additional simulation studies that vary the expressions for $\mathbf{E}(\mathbf{y}_L|\mathbf{x})$ and $\mathbf{E}(\mathbf{y}_U|\mathbf{x})$, as well as the distribution of ε_L (to be Beta(2,2) instead of Uniform(0,1)). Qualitatively the results are similar to what we report here.

We also perform a simulation study for the case that $d = 2$ with the following data generating process (DGP2):

$$\mathbf{Y} = \begin{bmatrix} 0.90 + 1.27\mathbf{x} + 10.18\mathbf{x}^2 \\ 0.60 - 1.00\mathbf{x} - 5.18\mathbf{x}^2 \end{bmatrix} + B_\xi,$$

where B_ξ is a ball of radius 1 centered at the random vector ξ , and ξ is uniformly distributed on the unit ball in \mathbb{R}^2 . As in the previous simulation, \mathbf{x} is drawn from a Beta distribution with support shifted to be $[-1, 1]$ and with shape parameters (2, 4). We let the sample size $n = 200, 500, 1000, 2000$. For values of $x_0 = 0, 0.2, 0.4, 0.6$ we evaluate the coverage probability of the error bands in equation (19) for $v = (1, 0)$, $v = (1, 1)/\sqrt{2}$, and $v = (0, 1)$. To conserve space, we report the results for $v = (1, 0)$ in Table 2 here, and for $v = (1, 1)/\sqrt{2}$ and $v = (0, 1)$, respectively, in Tables E.6 and E.7 in Appendix E. Overall the results are qualitatively similar to those reported for DGP1: once the bandwidth is undersmoothed and sample size is sufficiently large, coverage becomes valid.

7. Empirical Application

We demonstrate the usefulness of our approach with an empirical illustration that studies the association between cancer treatment outcomes and certain gene expression measures.

A key outcome of interest in cancer treatment research is the progression-free survival (PFS), which is defined as the time measured in months from baseline until tumor progression or death (whichever occurs first). Tumor progression is defined as an increase in the diameter of the tumor lesions of 20% compared with the smallest diameters of all previous tumor assessments or the appearance of

Table 2: Coverage probability at 95% nominal level using cross-validation for DGP2 with $v = (1, 0)$.

sample size	x_0	Coverage of $M(x_0)$	Coverage of $\mathbf{E}(\hat{M}(x_0))$	Coverage of $M(x_0)$ with $h = 1/2h_{n,CV}$	Coverage of $M(x_0)$ with $h = 1/3h_{n,CV}$
200	-0.4	0.8290	0.8395	0.8960	0.9725
	0	0.8515	0.8760	0.8525	0.9530
	0.2	0.9290	0.9360	0.9840	0.9985
	0.4	0.9085	0.9290	0.9220	0.9835
500	-0.4	0.8580	0.8665	0.9345	0.9805
	0	0.9195	0.9275	0.9745	0.9960
	0.2	0.9260	0.9325	0.9730	0.9945
	0.4	0.9210	0.9270	0.8965	0.9675
1000	-0.4	0.8830	0.8910	0.9315	0.9820
	0	0.9055	0.9125	0.9330	0.9785
	0.2	0.9210	0.9255	0.9425	0.9875
	0.4	0.9325	0.9345	0.9120	0.9725
2000	-0.4	0.8805	0.8835	0.9495	0.9875
	0	0.8900	0.8985	0.9355	0.9860
	0.2	0.9220	0.9300	0.9490	0.9915
	0.4	0.9270	0.9360	0.9595	0.9900

new lesions, as measured by CT-scans or MRIs (this is called RECIST criterion in the medical literature, see [32]). However, due to ethical and cost constraints, CT-scans and MRIs cannot be performed daily, but rather scheduled every 3 to 6 months. Hence, the PFS of patients can only be measured by intervals (with the true PFS falling between the last assessment without tumor progression and the assessment with progression), and no information is available on the distribution of true PFS within the interval. In contrast, the PFS of patients who died without tumor progression is measured exactly.

The question that we focus on in this paper is part of a subproject of the Swiss Cancer Research Group (SAKK) 19/09 for anti-cancer treatment regimens described in [33]. This subproject is concerned with finding, out of a total of 202 investigated genes, those whose baseline expression affects patient's PFS differently in two treatment arms described below. Genes expression is evaluated by isolating RNA from baseline tumor tissue sections and processing it for gene expression analysis using the Nanostring nCounter® System (Nanostring

Technologies), including 6 housekeeping genes.⁴ The gene expression measure that we report and use for our analysis is the \log_2 of the output of Nanostring.

It is worth mentioning that classical statistical methods of survival analysis, such as Cox regression or the accelerated failure time model, can also be applied to this data (and we do so below). These models are typically implemented with a parametric or semi-parametric specification of the hazard rate to construct the likelihood function. For example, the Cox proportional hazard model [34] assumes a hazard rate that is constant over time, and the resulting survival data follow a Markovian process; the accelerated failure time model posits an acceleration factor that is constant over time. The probability of censoring can then be calculated based on the functional form assumption. For example, the PFS variable in our example is usually treated as an interval censored data, for which one can construct the likelihood function and obtain point identified estimates of the model’s parameters, and then back out the implied conditional expectation of the treatment outcome given gene expression. In contrast, our method provides a consistent estimator of the set of admissible values for the conditional expectation of treatment outcome given gene expression, as well as $1 - \alpha$ pointwise confidence bands for it as in (16), without making any assumption on how PFS is distributed over the measured intervals that it is known to belong to, nor how it is related to the genes, as these assumptions may fail to hold in a given application.⁵

We use a novel dataset that follows 132 patients who were accrued between November 2010 and July 2014 to the SAKK 19/09 clinical trial for anti-cancer treatment regimens described in [33]. These patients are affected by advanced non-squamous non-small cell lung cancer and present an epidermal growth factor receptor (EGFR) of the wild type. Excluding 3 patients with protocol violations, 77 patients were treated with the drug Bevacizumab plus chemotherapy (C1) and 52 were treated with chemotherapy alone (C2). The question of interest

⁴See <https://www.nanostring.com> for a description of this method.

⁵[35] point out that individual heterogeneity and hazard rate cannot be jointly non-parametrically point identified.

Table 3: Descriptive statistics for interval-valued PFS and genes PTGS2 and CDC25A; \mathbf{y} denotes the progression-free survival (time from baseline until tumor progression or death), \mathbf{y}_L is last assessment without tumor progression, and \mathbf{y}_U is the assessment with tumor progression.

variable	mean	stdErr	max	min	N
\mathbf{y}_L	7.62	9.08	52.40	0	95
\mathbf{y}_U	9.25	9.65	55.16	0.23	95
CDC25A	7.23	2.76	14.22	0	95
PTGS2	8.66	1.90	13.37	2.86	95

of the SAKK 19/09 subproject that we revisit in this section is whether the
415 gene expression of PTGS2 (COX2) at baseline affects differently patient's PFS
in the two treatment arms. The gene PTGS2 (COX2) is frequently expressed
in lung cancer patients and the drug Bevacizumab directly interacts with the
COX2 pathway. One speculates that in patients with a high expression of
COX2 the tumor cells are predominately dependent on this signaling pathway
420 for proliferation and the use of Bevacizumab has a more pronounced effect.
Vice-versa, if COX2 is only expressed at a low level, this could reflect a tumor
that is not dependent on this inflammatory pathway and therefore the use of
Bevacizumab is not beneficial. Another gene of interest (whose effect on cancer
treatment efficacy has not been previously analyzed) is CDC25A, which is a
425 key regulator of the cells cycles. One speculates that overexpression of gene
CDC25A is associated with a poorer prognosis with regard to its biological role.

In our analysis, \mathbf{y} = PFS, \mathbf{y}_L is the time of the last assessment without
tumor progression, and \mathbf{y}_U is the time of the assessment with tumor progres-
sion. Table 3 reports descriptive statistics for these data. The sample used
for the analysis is constituted by 99 patients, from which four were excluded
because they were still alive at the last follow up (and therefore for these pa-
tients $\mathbf{y}_{iU} = \infty$). Of the sample used for our analysis, 58 patients were treated

following protocol C1, and 37 following protocol C2. Because durations are non-negative by definition while local linear regression smoothers may yield negative predictions, we work with the natural logarithm of our data, adjusted as follows

$$\tilde{\mathbf{y}}_k = \ln(\mathbf{y}_k + 0.033), \quad k = L, U$$

where we add 0.033 because for some individuals $\mathbf{y}_L = 0$. The choice of 0.033 is motivated by the unit of measure for \mathbf{y} , which is months: following the convention in the medical literature, we add one day (approximately 0.033 months).

430 The results of the analysis are reported in the top panels of Figure 2 for the gene PTGS2 (COX2), with panel A reporting the results using the Accelerated Failure Time (AFT) model, and Panel B reporting our set-valued local linear regression estimator. The bottom panels of Figure 2 report the results for the gene CDC25A, with panel C reporting the results using the AFT model, and
435 Panel D reporting our set-valued local linear regression estimator.

We first comment on the comparison between the standard AFT model and our set-valued estimator in terms of the shape of the predicted conditional PFS. For the PTGS2 (COX2) gene, the patterns are similar, although we uncover a more markedly nonlinear relation (especially for treatment C1). For the gene
440 CDC25A, the pattern uncovered by the AFT method and our method are similar for treatment C2, but for treatment C1 we uncover a remarkably more nonlinear relationship.

The results of the AFT analysis suggest that the use of Bevacizumab in cancer treatment is quite beneficial for patients with moderate to relatively high
445 (6-10) expression of gene PTGS2 (COX2), although the benefit seems to taper off at extremely high levels of the gene. Similarly, at medium to high levels (6-12) of expression of gene CDC25A the use of Bevacizumab seems beneficial, while at low levels of the gene the two treatment arm's effects are not significantly different. Our results, however, suggest that these findings might result from
450 the functional form assumptions: for the gene PTGS2 (COX2) we find that

for patients with moderate to relatively high (6-10) levels of the gene the set-valued estimates are consistent with a beneficial effect of Bevacizumab, but the confidence bands overlap, suggesting that the difference is not statistically significant. For the gene CDC25A we find that for CDC25A levels between 9
455 and 10, Bevacizumab is (statistical significantly) beneficial, but not at other levels of gene expression.

We note, however, that the results of this analysis are retrospective. To confirm the medical findings, a prospective randomized clinical trial needs to be carried out. We also highlight a drawback of our method: it is not able
460 to handle survival data censored on the right, where the observations become half-lines unbounded on the right. In our example such observations have been eliminated.

8. Conclusions

This paper has introduced local linear regression smoothing for set-valued
465 data. We have established consistency of the set-valued estimator, derived its mean squared error, and its (pointwise) asymptotic distribution. We have extended the cross-validation method for bandwidth selection to the case of set-valued local linear regression, and examined the finite sample properties of our estimator in a Monte Carlo exercise. We have illustrated the usefulness of our
470 method in an empirical illustration studying the effect of gene expression on cancer therapy outcomes.

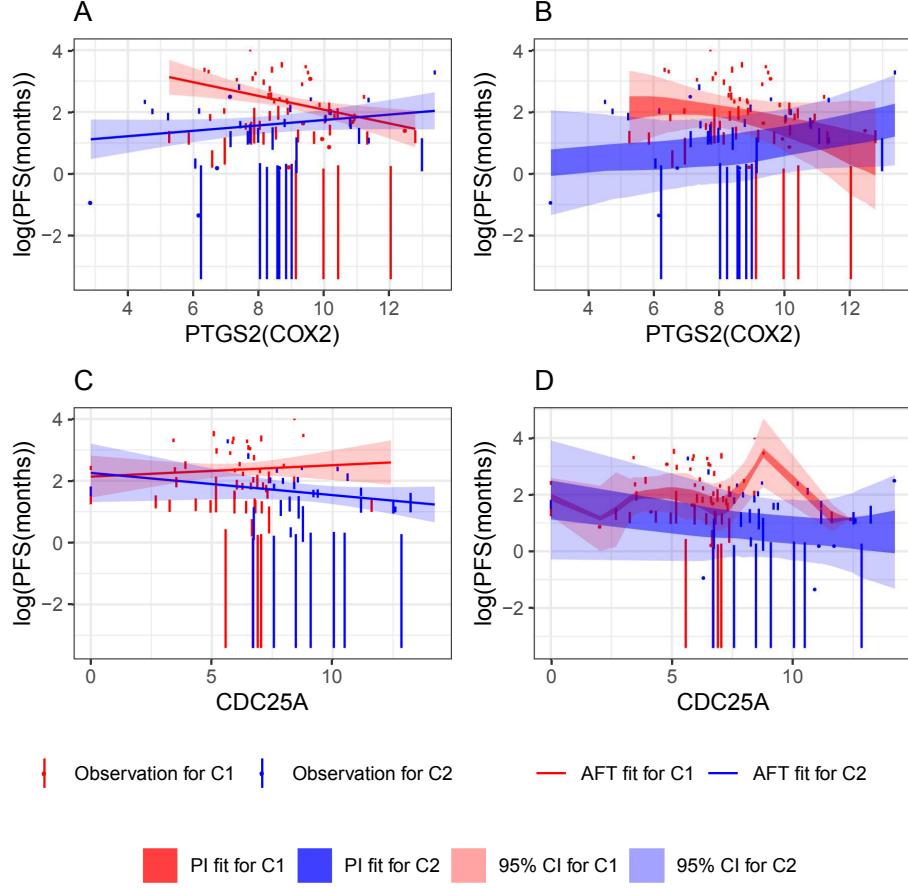


Figure 2: Results of the analysis for the genes PTGS2 and CDC25A (\log_2 of the Nanostring output)

Appendix A. Proofs of Main Results

Proof of Proposition 4.1. Our proof builds on [5, Eqs. (6.4), (6.6) and (6.13)]. Since the kernel is assumed to have a compact support, we have $\int z^{2r} K(z) dz < \infty$ for all $r \geq 0$. For any integer $r \geq 1$,

$$\mathbf{s}_j = \mathbf{E}\mathbf{s}_j + h_n^{j+1} \mathcal{O}_r(1/\sqrt{nh_n}), \quad j = 0, 1, 2, \quad (\text{A.1})$$

as $n \rightarrow \infty$, $h_n \rightarrow 0$ and $nh_n \rightarrow \infty$. The expectations of \mathbf{s}_j can be calculated as follows:

$$\begin{aligned}\mathbf{E}\mathbf{s}_0 &= h_n \int K(z) f(zh_n + x_0) dz = h_n \int K(z) (f(x_0) + \mathcal{O}(h_n)) dz = h_n [f(x_0) + \mathcal{O}(h_n)], \\ \mathbf{E}\mathbf{s}_1 &= h_n^2 \int zK(z) f(zh_n + x_0) dz = h_n^2 \int zK(z) (f(x_0) + \mathcal{O}(h_n)) dz = h_n^2 \mathcal{O}(h_n), \\ \mathbf{E}\mathbf{s}_2 &= h_n^3 \int z^2 K(z) f(zh_n + x_0) dz = h_n^3 \int z^2 K(z) (f(x_0) + \mathcal{O}(h_n)) dz = h_n^3 (f(x_0) \text{Var}_K + \mathcal{O}(h_n)).\end{aligned}$$

In view of (A.1), for an integer $r \geq 1$,

$$\mathbf{s}_j = h_n^{j+1} \left(f(x_0) \int z^j K(z) dz + \mathcal{O}_r(h_n + \frac{1}{\sqrt{nh_n}}) \right), \quad j = 0, 1, 2. \quad (\text{A.2})$$

Thus,

$$\mathbf{s}_0 = h_n f(x_0) (1 + \mathcal{O}_r(1)), \quad (\text{A.3})$$

$$\mathbf{s}_1 = h_n^2 \mathcal{O}_r(1), \quad (\text{A.4})$$

$$\mathbf{s}_2 = h_n^3 f(x_0) \text{Var}_K (1 + \mathcal{O}_r(1)). \quad (\text{A.5})$$

It is easy to see that

$$\sum_{i=1}^n \ell_i = \frac{\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2}{\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4}}.$$

Moreover, for a sufficiently large r ,

$$\frac{h_n^4}{\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4}} = \frac{1}{f(x_0)^2 \text{Var}_K} + \mathcal{O}_r(1), \quad (\text{A.6})$$

cf. [5, Eq. (6.6)]. In view of (A.3), (A.4), and (A.5),

$$\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 = h_n^4 f(x_0)^2 \text{Var}_K (1 + \mathcal{O}_r(1)). \quad (\text{A.7})$$

By (11),

$$\sum_{i=1}^n \ell_i^2 = \frac{\sum_{i=1}^n \kappa_{in}^2 (\mathbf{s}_2 - (\mathbf{x}_i - x_0) \mathbf{s}_1)^2}{n^2 (\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4})^2} = \frac{\mathbf{s}_2^2 \mathbf{s}_0^*}{n (\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4})^2} + \frac{(-2 \mathbf{s}_2 \mathbf{s}_1 \mathbf{s}_1^* + \mathbf{s}_1^2 \mathbf{s}_2^*)}{n (\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4})^2}, \quad (\text{A.8})$$

where

$$\mathbf{s}_j^* = \frac{1}{n} \sum_{i=1}^n \kappa_{in}^2 (x_i - x_0)^j = h_n^{j+1} \left(f(x_0) \int z^j K^2(z) dz + \mathcal{O}_r(1) \right), \quad j = 0, 1, 2.$$

Furthermore, (A.2) implies that

$$\mathbf{s}_2^2 \mathbf{s}_0^* = h_n^7 f^3(x_0) (\text{Var}_K)^2 \int K^2(z) dz + h_n^7 \mathcal{O}_{r/2}(1).$$

Combining this with (A.6) and letting $r = 4$, we obtain

$$\begin{aligned} \mathbf{E} \left(\frac{\mathbf{s}_2^2 \mathbf{s}_0^*}{n (\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4})^2} \right) &= \frac{h_n^7 f^3(x_0) (\text{Var}_K)^2 \int K^2(z) dz}{n h_n^8 f^4(x_0) (\text{Var}_K)^2} + \frac{h_n^7}{n h_n^8} \mathcal{O}(1) \\ &= \frac{\int K^2(z) dz}{n h_n f(x_0)} + \mathcal{O} \left(\frac{1}{n h_n} \right). \end{aligned}$$

Since $\int z K(z) dz = 0$,

$$-2 \mathbf{s}_2 \mathbf{s}_1 \mathbf{s}_1^* = h_n^7 (f(x_0) \text{Var}_K + \mathcal{O}_8(1)) \mathcal{O}_8(1) (f(x_0) \int z^j K^2(z) dz + \mathcal{O}_4(1)) = h_n^7 \mathcal{O}_2(1).$$

Analogously, $\mathbf{s}_1^2 \mathbf{s}_2^* = h_n^7 \mathcal{O}_2(1)$. Both these terms are as small as the minor term of $\mathbf{s}_2^2 \mathbf{s}_0^*$. Therefore, (A.8) is dominated by its first term, whence (12) holds. \square

Proof of Proposition 4.2. By definition, $\ell_i < 0$ if and only if $\mathbf{s}_2 - (\mathbf{x}_i - x_0) \mathbf{s}_1 < 0$.

Hence,

$$\begin{aligned}
\mathbf{E}\left(\sum_{i=1}^n \ell_i^-\right)^2 &= \mathbf{E}\left(\sum_{i=1}^n -\ell_i \mathbf{1}\{\mathbf{s}_2 - (\mathbf{x}_i - x_0)\mathbf{s}_1 < 0\}\right)^2 \leq n\mathbf{E}\left(\sum_{i=1}^n \ell_i^2 \mathbf{1}\{\mathbf{s}_2 - (\mathbf{x}_i - x_0)\mathbf{s}_1 < 0\}\right) \\
&\leq n\mathbf{E}\left(\sum_{i=1}^n \ell_i^2 \mathbf{1}\{\mathbf{s}_2 < c_K h_n |\mathbf{s}_1|\}\right) = n\mathbf{E}\left(\mathbf{1}\{\mathbf{s}_2 < c_K h_n |\mathbf{s}_1|\} \sum_{i=1}^n \ell_i^2\right) \\
&\leq n\sqrt{\mathbf{P}(\mathbf{s}_2 < c_K h_n |\mathbf{s}_1|)} \left(\mathbf{E}\left(\sum_{i=1}^n \ell_i^2\right)^2\right)^{1/2}, \tag{A.9}
\end{aligned}$$

where the second inequality relies on Assumption A and the last one follows from the Chebyshev inequality. Using (A.2), we have, for an integer $r \geq 1$,

$$\begin{aligned}
\mathbf{s}_1 &= h_n^2 \left(\mathcal{O}(h_n) + \mathcal{O}_r(1/\sqrt{nh_n}) \right), \\
\mathbf{s}_2 &= h_n^3 \left(f(x_0) \text{Var}_K + \mathcal{O}(h_n) + \mathcal{O}_r(1/\sqrt{nh_n}) \right).
\end{aligned}$$

Hence,

$$\mathbf{P}(\mathbf{s}_2 < c_K h_n |\mathbf{s}_1|) \tag{A.10}$$

$$\begin{aligned}
&\leq \mathbf{P}\left(f(x_0) \text{Var}_K + \mathcal{O}(h_n) + \mathcal{O}_r(1/\sqrt{nh_n}) < |\mathcal{O}(h_n)| + \left|\mathcal{O}_r(1/\sqrt{nh_n})\right|\right) \\
&= \mathbf{P}\left(f(x_0) \text{Var}_K < |\mathcal{O}(h_n)| + \left|\mathcal{O}_r(1/\sqrt{nh_n})\right|\right). \tag{A.11}
\end{aligned}$$

For sufficiently large n , there exist a ξ with $0 < \xi < f(x_0) \text{Var}_K$ so that $|\mathcal{O}(h_n)| \leq \xi$ for all sufficiently large n . Building on (A.11), the Markov inequality and the definition of $\mathcal{O}_r(a_n)$ yield that

$$\begin{aligned}
\mathbf{P}(\mathbf{s}_2 < c_K h_n |\mathbf{s}_1|) &\leq \mathbf{P}\left(f(x_0) \text{Var}_K < \xi + \left|\mathcal{O}_r(1/\sqrt{nh_n})\right|\right) \\
&= \mathbf{P}\left(\left|\mathcal{O}_r(1/\sqrt{nh_n})\right| > f(x_0) \text{Var}_K - \xi\right) \\
&\leq \frac{\sup_{f \in \mathcal{H}(1, \gamma)} \mathbf{E} \left|\mathcal{O}_r(1/\sqrt{nh_n})\right|^r}{(f(x_0) \text{Var}_K - \xi)^r} = \frac{c_r (1/\sqrt{nh_n})^r}{(f(x_0) \text{Var}_K - \xi)^r}
\end{aligned}$$

for a positive constant c_r . Therefore,

$$\mathbf{P}(\mathbf{s}_2 < c_K h_n |\mathbf{s}_1|) = \mathcal{O}\left((1/\sqrt{nh_n})^r\right). \quad (\text{A.12})$$

From the proof of Proposition 4.1 with $r = 8$, squaring and taking expectation,

$$\mathbf{E}\left(\sum_{i=1}^n \ell_i^2\right)^2 = \frac{1}{n^2 h_n^2} \left(\int K^2(z) dz\right)^2 (1 + \mathcal{O}(1)). \quad (\text{A.13})$$

Substituting (A.12) and (A.13) into (A.9),

$$\mathbf{E}\left(\sum_{i=1}^n \ell_i^-\right)^2 \leq \frac{1}{h_n} \int K^2(z) dz \sqrt{1 + \mathcal{O}(1)} \mathcal{O}\left((1/\sqrt{nh_n})^r\right),$$

475 which converges to 0 by choosing a sufficiently large r . □

Proof of Theorem 4.3. The squared bias can be written as

$$b_{x_0}^2(v) = \mathbf{E}[(b_1 + b_2)^2],$$

for $b_1 = \sum_{i=1}^n \ell_i(s(M(\mathbf{x}_i), v) - s(M(x_0), v))$ and $b_2 = \sum_{i=1}^n \ell_i^- w(M(\mathbf{x}_i), v)$. We have

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \kappa_{in}(\mathbf{s}_2 - (\mathbf{x}_i - x_0)\mathbf{s}_1)(s(M(\mathbf{x}_i), v) - s(M(x_0), v)) \\ &= \frac{1}{n} \sum_{i=1}^n \kappa_{in}(\mathbf{s}_2 - (\mathbf{x}_i - x_0)\mathbf{s}_1)(s(M(\mathbf{x}_i), v) - s(M(x_0), v) + s'(M(x_0), v)(\mathbf{x}_i - x_0)) \\ &= h_n^6 f(x_0) \text{Var}_K a_n + \mathcal{O}_4(h_n^6), \end{aligned}$$

where

$$a_n = h_n^{-3} \mathbf{E} \left(s(M(\mathbf{x}), v) - s(M(x_0), v) - s'(M(x_0), v)(\mathbf{x} - x_0) K\left(\frac{\mathbf{x} - x_0}{h_n}\right) \right).$$

By (A.6), and using the definition of \mathcal{O}_r , we have

$$\mathbf{E}b_1^2 = \mathbf{E} \left(\frac{\frac{1}{n} \sum_{i=1}^n \kappa_{in}(\mathbf{s}_2 - (\mathbf{x}_i - x_0)\mathbf{s}_1)(m_v(\mathbf{x}_i) - m_v(x_0))}{\mathbf{s}_2\mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4}} \right)^2 = \left(\frac{U_n}{f(x_0)} \right)^2 h_n^4 + \mathcal{O}(h_n^4),$$

where, taking a Taylor expansion,

$$U_n = h_n^{-2} \left(\frac{1}{2} s''(M(x_0), v) \text{Var}_K f(x_0) h_n^2 + \mathcal{O}(h_n^2) \right).$$

Therefore,

$$\mathbf{E}b_1^2 = \frac{1}{4} s''(M(x_0), v)^2 (\text{Var}_K)^2 h_n^4 + \mathcal{O}(h_n^4), \quad (\text{A.14})$$

cf. the proof of [5, Theorem 3].

By Proposition 4.2,

$$\mathbf{E}b_2^2 \leq w_{\max}^2 \mathbf{E} \left(\sum_{i=1}^n \ell_i^- \right)^2 = \frac{1}{h_n} \mathcal{O} \left((1/\sqrt{nh_n})^r \right), \quad (\text{A.15})$$

where w_{\max} is a finite deterministic bound on the width of $M(\mathbf{x})$ in any direction $v \in \mathbb{S}^{d-1}$ resulting from Assumption B. By the Cauchy-Schwarz inequality, (A.15) and (A.14),

$$\mathbf{E}(b_1 b_2) \leq \sqrt{\mathbf{E}b_1^2 \mathbf{E}b_2^2} = \frac{1}{2} \left(s''(M(x_0), v)^2 (\text{Var}_K)^2 h_n^4 + \mathcal{O}(h_n^4) \right)^{1/2} h_n^{-1/2} \mathcal{O} \left((1/\sqrt{nh_n})^{r/2} \right),$$

which, for sufficiently large r and given that $h_n = cn^{-\beta}$, is of a smaller order than h_n^4 . Thus,

$$\int_{\mathbb{S}^{d-1}} b_{x_0}^2(v) dv = \frac{1}{4} \int_{\mathbb{S}^{d-1}} s''(M(x_0), v)^2 dv (\text{Var}_K)^2 h_n^4 + \mathcal{O} \left(h_n^4 + \frac{1}{nh_n} \right). \quad (\text{A.16})$$

Now we bound the variance of the estimator splitting (10) into the sum of three terms. By Proposition 4.1, the first term is

$$\mathbf{E} \left(\sum_{i=1}^n \ell_i \varepsilon_i(v) \right)^2 = \mathbf{E} \sum_{i=1}^n \ell_i^2 C(v, v) = \frac{1}{nh_n f(x_0)} C(v, v) \int K^2(z) dz + \mathcal{O} \left(\frac{1}{nh_n} \right).$$

The second term is

$$\mathbf{E} \sum_{1 \leq i < j \leq n} \ell_i \ell_j^- \varepsilon_i(v) (\varepsilon_j(v) + \varepsilon_j(-v)) = 0.$$

Finally, consider

$$\begin{aligned} \mathbf{E} \left(\sum_{i=1}^n \ell_i^- (\varepsilon_i(v) + \varepsilon_i(-v)) \right)^2 &= (C(v, v) + 2C(v, -v) + C(-v, -v)) \mathbf{E} \sum_{i=1}^n (\ell_i^-)^2 \\ &\leq 4\sigma_{\max}^2 \mathbf{E} \sum_{i=1}^n (\ell_i^-)^2 \leq 4\sigma_{\max}^2 \mathbf{E} \left(\sum_{i=1}^n \ell_i^- \right)^2 \\ &= 4\sigma_{\max}^2 h_n^{-1} \mathcal{O} \left((1/\sqrt{nh_n})^r \right). \end{aligned}$$

For a large r , $(nh_n)^{(-r/2)}$ is of smaller order than $(nh_n)^{-1}$. Hence,

$$\int_{\mathbb{S}^{d-1}} \sigma_{x_0}^2(v) dv = \frac{1}{nh_n f(x_0)} \int_{\mathbb{S}^{d-1}} C(v, v) dv \int K^2(z) dz + \mathcal{O} \left(\frac{1}{nh_n} \right),$$

and the result follows by adding (A.16) to it. \square

Proof of Theorem 4.4. It suffices to establish the convergence of one-dimensional distributions; the weak convergence of finite dimensional distributions follows from the Cramér–Wold device, and the functional convergence is established by

480

bounding the Lipschitz constants of the processes as in [4, Theorem 3.2.1].

First, decompose

$$\begin{aligned} s(\hat{M}, v) - s(M(x_0), v) &= \sum_{i=1}^n \ell_i s(\mathbf{Y}_i, v) + \sum_{i=1}^n \ell_i^- w(\mathbf{Y}_i, v) - s(M(x_0), v) \\ &= \sum_{i=1}^n \ell_i s(M(\mathbf{x}_i), v) + \sum_{i=1}^n \ell_i \varepsilon_i(v) + \sum_{i=1}^n \ell_i^- w(\mathbf{Y}_i, v) - s(M(x_0), v). \end{aligned} \tag{A.17}$$

By Proposition 4.2, noticing that the L_2 -convergence implies the convergence

in probability, and choosing r large enough, we have that

$$\sum_{i=1}^n \ell_i^- w(\mathbf{Y}_i, v) \leq w_{\max} \sum_{i=1}^n \ell_i^- = \mathcal{O}_p(1/\sqrt{nh_n}).$$

Using a Taylor expansion,

$$s(M(\mathbf{x}_i), v) = s(M(x_0), v) + (\mathbf{x}_i - x_0)s'(M(x_0), v) + \frac{1}{2}(\mathbf{x}_i - x_0)^2 s''(M(x_0), v) + R(x_0, \mathbf{x}_i, v),$$

where the remainder term $R(x_0, \mathbf{x}_i, v)$ is of a smaller order than $\frac{1}{2}(\mathbf{x}_i - x_0)^2 s''(M(x_0), v)$.

Since the local linear estimator satisfies $\sum_{i=1}^n \ell_i(\mathbf{x}_i - x_0) = 0$, we have

$$\begin{aligned} & \sum_{i=1}^n \ell_i s(M(\mathbf{x}_i), v) + \sum_{i=1}^n \ell_i \varepsilon_i(v) - s(M(x_0), v) \\ &= \sum_{i=1}^n \ell_i (s(M(\mathbf{x}_i), v) - s(M(x_0), v)) - \frac{n^{-4}}{\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4}} s(M(x_0), v) + \sum_{i=1}^n \ell_i \varepsilon_i(v) \\ &= \sum_{i=1}^n \ell_i \left(\frac{1}{2}(\mathbf{x}_i - x_0)^2 s''(M(x_0), v) + R(x_0, \mathbf{x}_i, v) + \varepsilon_i(v) \right) - \frac{n^{-4}}{\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4}} s(M(x_0), v). \end{aligned}$$

Since for a sequence of $\{Z_n, n \geq 1\}$ of square integrable random variables

$$Z_n = \mathbf{E}Z_n + \mathcal{O}_p(\sqrt{\text{Var } Z_n}),$$

(A.2) yields that

$$\mathbf{s}_j = h_n^{j+1} f(x_0) \int z^j K(z) dz (1 + \mathcal{O}_p(1)), \quad j = 0, 1, 2, 3. \quad (\text{A.18})$$

By (A.7) and since $nh_n \rightarrow \infty$, we have

$$\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4} = h_n^4 \text{Var}_K f^2(x_0) (1 + \mathcal{O}_p(1)). \quad (\text{A.19})$$

Therefore,

$$\frac{n^{-4}}{\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4}} s(M(x_0), v) = \mathcal{O}_p(n^{-4} h_n^{-4}) = \mathcal{O}_p(n^{-3} h_n^{-3}).$$

Combining (A.18) and (A.19), we have

$$\begin{aligned}
& \sum_{i=1}^n \ell_i \left(\frac{1}{2} (\mathbf{x}_i - x_0)^2 s''(M(x_0), v) + R(x_0, \mathbf{x}_i, v) + \varepsilon_i(v) \right) \\
&= \left(\frac{1}{2} (s_2^2 - \mathbf{s}_3 \mathbf{s}_1) s''(M(x_0), v) + \frac{1}{n} \sum_{i=1}^n \kappa_{in} (\mathbf{s}_2 - (\mathbf{x}_i - x_0) \mathbf{s}_1) \varepsilon_i(v) \right) (\mathbf{s}_2 \mathbf{s}_0 - \mathbf{s}_1^2 + n^{-4})^{-1} \\
&= \frac{1}{2} \text{Var}_K s''(M(x_0), v) h_n^2 (1 + \mathcal{O}_p(1)) + \frac{1}{n h_n f(x_0)} \sum_{i=1}^n \kappa_{in} \varepsilon_i(v) (1 + \mathcal{O}_p(1)).
\end{aligned} \tag{A.20}$$

By the central limit theorem,

$$\frac{1}{\sqrt{n h_n}} \sum_{i=1}^n \kappa_{in} \varepsilon_i \tag{A.21}$$

converges in distribution to the centered normal random variable with variance equal to that of $\zeta(v)$. The combination of (A.17), (A.19), (A.20) and (A.21) yields the result. \square

485 Appendix B. Deterministic design points

When the design points $\mathbf{x}_i = x_i$, $i = 1, \dots, n$, are deterministic⁶, (9) turns into

$$b_{x_0}^2(v) = \left(\sum_{i=1}^n \ell_i (s(M(x_i), v) - s(M(x_0), v)) + \sum_{i=1}^n \ell_i^- w(M(x_i), v) \right)^2.$$

Since $K(\cdot)$ has compact support in $[-c_K, c_K]$, we have $\ell_i = 0$ if $|x_i - x_0| > c_K h_n$. It is easy to see that all weights are nonnegative if and only if

$$\sum \kappa_{in} \left(\frac{x_i - x_0}{h_n} \right)^2 \geq \left| \sum \kappa_{in} \frac{x_i - x_0}{h_n} \right|.$$

This assumption means that the sample rescaled around each point to lie in

⁶Because with deterministic design $\mathbf{x}_i = x_i$, $i = 1, \dots, n$, \mathbf{s}_j , $j = 0, 1, 2$ and κ_{in} , $i = 1, \dots, n$ are also deterministic and we write $\mathbf{s}_j = s_j$ and $\kappa_{in} = \kappa_{in}$.

the range $[-1, 1]$ has the variance that dominates the absolute value of the expectation. For this, the rescaled points should be sufficiently balanced on the left and on the right of x_0 . The assumption can be alternatively expressed as

$$\frac{s_2}{h_n^3} \geq c_K \left| \frac{s_1}{h_n^2} \right|.$$

It holds when $s_1/h_n^2 \rightarrow 0$ as $n \rightarrow \infty$.

By a direct computation, it is possible to show that, in the regular design case, the weights are nonnegative for all n .

Proposition Appendix B.1. *Consider the local linear setting with uniform kernel supported on $[-c_K, c_K]$ and equally spaced (regular) design points x_1, \dots, x_n on a bounded interval I . If $1/n \leq c_K h_n \leq 1$, then $\ell_i(x_0) \geq 0$ for all i , n and each*

$$x_0 \in I_n = \{x \in I : [x - c_K h_n, x + c_K h_n] \subset I\}.$$

In case of deterministic design points in a bounded interval I , the following
490 assumptions are often imposed; they appear as (LP1)-(LP2) in [3].

Assumption E (Design points). *The design points x_1, \dots, x_n are such that:*

- (i) *There exists $\lambda_0 > 0$ such that all eigenvalues of $\mathcal{B}_{n x_0}$ are greater than or equal to λ_0 for all sufficiently large n and all $x_0 \in I$.*
- (ii) *There exists $a_0 > 0$ such that, for any interval $J \subset I$ and all $n > 1$,*

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i \in J} \leq a_0 \max(\text{Leb}(J)/\text{Leb}(I), 1/n),$$

where $\text{Leb}(\cdot)$ denotes the Lebesgue measure.

495 We impose the following assumption on the response function.

Assumption F (Theoretical response function). *The function $M(x)$, $x \in I$, is defined on a bounded closed interval $I \subset \mathbb{R}$, and there exists $\gamma > 0$ such that,*

for all $v \in \mathbb{S}^{d-1}$, the derivative of $s(M(x), v)$ with respect to x is Lipschitz with constant γ .

500 The following result is similar to [3, Prop. 1.13] in the singleton-valued data framework.

Proposition Appendix B.2. *If $x_0 \in I_n$, $\ell_i \geq 0$ for all i , and Assumptions A, B, E and F are satisfied, then*

$$|b_{x_0}(v)| \leq c_K^2 C_* \gamma h_n^2, \quad \sigma_{x_0}^2(v) \leq \frac{\sigma_{\max}^2 C_*^2}{n h_n}$$

for sufficiently large n and $h_n \geq 1/(2n)$.

Proposition Appendix B.2 implies

$$\text{MSE}(x_0) \leq c_K^4 C_*^2 \gamma^2 h_n^4 + \frac{\sigma_{\max}^2 C_*^2}{n h_n}.$$

Therefore, the upper bound is minimized for the bandwidth given by

$$h_n^* = \left(\frac{\sigma_{\max}^2}{4c_K^4 \gamma^2} \right)^{\frac{1}{5}} n^{-\frac{1}{5}},$$

and the following result holds.

Theorem Appendix B.3. *If the bandwidth is chosen to be $h_n = \alpha n^{-\frac{1}{5}}$ for $\alpha > 0$ and Assumptions A, B, E hold, then*

$$\limsup_{n \rightarrow \infty} \sup_{x_0 \in I_n} \mathbf{E}[n^{\frac{2}{5}} L(\hat{M}(x), M(x))] \leq C_1 < \infty,$$

505 uniformly over all response functions satisfying Assumption F, where L is the loss function given by (6), C_1 is a constant depending only on γ , a_0 , λ_0 , σ_{\max}^2 , K_{\max} and α .

Appendix C. Local constant regression

In the local constant case, the weights $\ell_i = \kappa_{in}/(n s_0)$ are always nonnegative. Then the estimator $\hat{M}(x_0)$ can be constructed as the convex set whose

510 support functions is obtained by calculating the Nadaraya–Watson estimator for
the sample $s(\mathbf{Y}_i, v)$, $i = 1, \dots, n$, in each particular direction v . In other words,
 $\hat{\mathbf{M}}(x_0)$ is the sum of the observed sets \mathbf{Y}_i multiplied by nonnegative coefficients
 ℓ_i . Therefore, the bias and variance of the set-valued local constant estimator
can be obtained similarly to the singleton-valued data case. For this, it suffices
515 to assume that the function $s(M(x), v)$ is Lipschitz in x with the same constant
for all v , which is equivalent to requiring that $M(x)$, $x \in I$, is Lipschitz in the
Hausdorff metric.

Appendix D. Basic definitions from random set theory

A *random compact set* \mathbf{Y} is a map from $(\Omega, \mathfrak{F}, \mathbf{P})$ to $\mathcal{K}(\mathbb{R}^d)$ such that

$$\{\omega : \mathbf{Y}(\omega) \cap K \neq \emptyset\} \in \mathfrak{F}, \quad (\text{D.1})$$

for each compact set $K \subset \mathbb{R}^d$.

Random sets $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ are said to be independent and identically distributed if

$$\mathbf{P}(\mathbf{Y}_1 \cap K_1 \neq \emptyset, \dots, \mathbf{Y}_n \cap K_n \neq \emptyset) = \prod_{i=1}^n \mathbf{P}(\mathbf{Y}_i \cap K_i \neq \emptyset),$$

520 for all $K_1, \dots, K_n \in \mathcal{K}(\mathbb{R}^d)$ and $\mathbf{P}(\mathbf{Y}_i \cap K \neq \emptyset) = \mathbf{P}(\mathbf{Y}_j \cap K \neq \emptyset)$ for all
 $i \neq j \in \{1, \dots, n\}$ and $K \in \mathcal{K}(\mathbb{R}^d)$.

We define the Minkowski sum of two compact sets A_1 and A_2 in \mathbb{R}^d elementwise as

$$A + B = \{x + y : x \in A, y \in B\}.$$

We let $cA = \{cx : x \in A\}$ denote the scaling of A by $c \in \mathbb{R}$. Given a compact convex set (a *convex body*) $A \subset \mathbb{R}^d$, the *support function* of A is

$$s(A, v) = \sup_{a \in A} v^\top a, \quad v \in \mathbb{R}^d,$$

where $v^\top a$ denotes the scalar product. If A is convex, its support function

uniquely identifies A , because

$$A = \bigcap_{v \in \mathbb{S}^{d-1}} \{a \in \mathbb{R}^d : v^\top a \leq s(A, v)\}. \quad (\text{D.2})$$

Because $s(tA, v) = ts(A, v)$ for $t \geq 0$, the support function is often restricted to $v \in \mathbb{S}^{d-1}$. Note that

$$s(A_1 + A_2, v) = s(A_1, v) + s(A_2, v).$$

The width function of A is defined by

$$w(A, v) = s(A, v) + s(A, -v) = w(A, -v), \quad v \in \mathbb{S}^{d-1},$$

and it is easy to see that the width function is nonnegative. If $d = 1$, then A is a closed interval in \mathbb{R} , and the unit sphere $\mathbb{S}^{d-1} = \{-1, 1\}$ consists of two points.

In this case, the width function is the length of the interval.

525 A *random convex compact set* \mathbf{Y} is a map from $(\Omega, \mathfrak{F}, \mathbf{P})$ to $\mathcal{K}_{\mathcal{C}}(\mathbb{R}^d)$ satisfying equation (D.1). Its measurability is equivalent to the fact that $s(\mathbf{Y}, v)$ is a random variable for each $v \in \mathbb{R}^d$.

Appendix E. Additional simulation results

Table E.4: Coverage probability at 95% nominal level using cross-validation for a modified DGP1 with $\mathbf{y}_L = 0.90 + 1.27\mathbf{x} + 10.18\mathbf{x}^2 - \epsilon_L$, $\mathbf{y}_U = 0.90 + 1.27\mathbf{x} + 10.18\mathbf{x}^2 + \epsilon_U$, and $\epsilon_L, \epsilon_U \sim^{i.i.d.} Uniform[0, 1]$.

sample size	x_0	Coverage of $M(x_0)$	Coverage of $\mathbf{E}(\hat{M}(x_0))$	Coverage of $M(x_0)$ with $h = 1/2h_{n,CV}$	Coverage of $M(x_0)$ with $h = 1/3h_{n,CV}$
200	-0.4	0.8630	0.8540	0.9165	0.9690
	0	0.8965	0.8865	0.8790	0.9520
	0.2	0.9465	0.9405	0.9825	0.9980
	0.4	0.9330	0.9215	0.9200	0.9745
500	-0.4	0.8705	0.8595	0.9290	0.9755
	0	0.9460	0.9410	0.9760	0.9935
	0.2	0.9315	0.9280	0.9655	0.9910
	0.4	0.9415	0.9320	0.9260	0.9800
1000	-0.4	0.9070	0.9040	0.9525	0.9855
	0	0.8990	0.8985	0.9175	0.9695
	0.2	0.9205	0.9160	0.9425	0.9760
	0.4	0.8965	0.8940	0.9090	0.9570
2000	-0.4	0.8970	0.8925	0.9440	0.9820
	0	0.9305	0.9290	0.9585	0.9865
	0.2	0.9230	0.9215	0.9425	0.9815
	0.4	0.8925	0.8935	0.9040	0.9600

Table E.5: Coverage probability at 95% nominal level using cross-validation for a modified DGP1 with $\mathbf{y}_L = 0.90 + 1.27\mathbf{x} - \epsilon_L$, $\mathbf{y}_U = 0.90 + 1.27\mathbf{x} + 10.18\mathbf{x}^2 + \epsilon_U$, $\epsilon_L \sim \text{Beta}(2, 2)$ and $\epsilon_U \sim \text{Uniform}(0, 1)$.

sample size	x_0	Coverage of $M(x_0)$	Coverage of $\mathbf{E}(\hat{M}(x_0))$	Coverage of $M(x_0)$ with $h = 1/2h_{n,CV}$	Coverage of $M(x_0)$ with $h = 1/3h_{n,CV}$
200	-0.4	0.6510	0.8945	0.9515	0.9865
	0	0.6610	0.9125	0.9050	0.9770
	0.2	0.7495	0.9600	0.9920	0.9995
	0.4	0.7210	0.9565	0.9680	0.9960
500	-0.4	0.6255	0.8945	0.9575	0.9875
	0	0.7200	0.9445	0.9870	0.9995
	0.2	0.7355	0.9605	0.9825	0.9985
	0.4	0.7155	0.9575	0.9525	0.9880
1000	-0.4	0.6345	0.9175	0.9660	0.9955
	0	0.6485	0.9330	0.9625	0.9895
	0.2	0.6960	0.9580	0.9715	0.9945
	0.4	0.7025	0.9535	0.9545	0.9870
2000	-0.4	0.6195	0.9255	0.9710	0.9935
	0	0.6290	0.9360	0.9610	0.9905
	0.2	0.6605	0.9500	0.9750	0.9935
	0.4	0.6755	0.9600	0.9785	0.9955

Table E.6: Coverage probability at 95% nominal level using cross-validation for DGP2 with $\mathbf{v} = (1, 1)/\sqrt{2}$.

sample size	x_0	Coverage of $M(x_0)$	Coverage of $\mathbf{E}(\hat{M}(x_0))$	Coverage of $M(x_0)$ with $h = 1/2h_{n,CV}$	Coverage of $M(x_0)$ with $h = 1/3h_{n,CV}$
200	-0.4	0.8225	0.9475	0.9490	0.9870
	0	0.8150	0.9370	0.9400	0.9820
	0.2	0.7825	0.9170	0.9330	0.9835
	0.4	0.7310	0.9020	0.9265	0.9815
500	-0.4	0.8445	0.9495	0.9635	0.9890
	0	0.7655	0.9195	0.9525	0.9895
	0.2	0.7385	0.9150	0.9410	0.9830
	0.4	0.6820	0.8745	0.9475	0.9890
1000	-0.4	0.8230	0.9500	0.9595	0.9895
	0	0.7945	0.9350	0.9455	0.9825
	0.2	0.7270	0.9185	0.9580	0.9900
	0.4	0.6830	0.8645	0.9290	0.9825
2000	-0.4	0.7965	0.9440	0.9480	0.9900
	0	0.7925	0.9430	0.9390	0.9860
	0.2	0.7485	0.9370	0.9390	0.9845
	0.4	0.7370	0.9250	0.9515	0.9890

Table E.7: Coverage probability at 95% nominal level using cross-validation for DGP2 with $v = (0, 1)$.

sample size	x_0	Coverage of $M(x_0)$	Coverage of $\mathbf{E}(\hat{M}(x_0))$	Coverage of $M(x_0)$ with $h = 1/2h_{n,CV}$	Coverage of $M(x_0)$ with $h = 1/3h_{n,CV}$
200	-0.4	0.8395	0.9450	0.9485	0.9875
	0	0.8085	0.9160	0.9230	0.9765
	0.2	0.7815	0.9090	0.9445	0.9840
	0.4	0.7405	0.8945	0.9310	0.9820
500	-0.4	0.8020	0.9395	0.9530	0.9875
	0	0.7995	0.9330	0.9545	0.9905
	0.2	0.7550	0.9210	0.9380	0.9805
	0.4	0.7215	0.9025	0.9495	0.9875
1000	-0.4	0.8175	0.9485	0.9560	0.9905
	0	0.7900	0.9405	0.9420	0.9870
	0.2	0.7290	0.9345	0.9535	0.9865
	0.4	0.7070	0.8830	0.9415	0.9895
2000	-0.4	0.7945	0.9440	0.9475	0.9895
	0	0.7935	0.9430	0.9395	0.9860
	0.2	0.7495	0.9375	0.9400	0.9845
	0.4	0.7355	0.9245	0.9515	0.9890

References

- [1] F. T. Juster, R. Suzman, An overview of the health and retirement study, Journal of Human Resources 30 (Supplement) (1995) S7–S56.
- [2] C. F. Manski, Partial Identification of Probability Distributions, Springer, New York, 2003.
- [3] A. B. Tsybakov, Introduction to Nonparametric Estimation, Springer, New York, 2009.
- [4] I. Molchanov, Theory of Random Sets, 2nd Edition, Springer, London, 2017.
- [5] J. Fan, Local linear regression smoothers and their minimax efficiencies, Ann. Statist. 21 (1993) 196–216.
URL <http://dx.doi.org/10.1214/aos/1176349022>
- [6] J. Fan, I. Gijbels, Local Polynomial Modelling and Its Applications, Chapman & Hall, London, 1996.

- [7] J. Fan, I. Gijbels, Variable bandwidth and local linear regression smoothers, *Ann. Statist.* 20 (1992) 2008–2036.
 545 URL <http://dx.doi.org/10.1214/aos/1176348900>
- [8] A. Beresteanu, F. Molinari, Asymptotic properties for a class of partially identified models, *Econometrica* 76 (2008) 763–814.
- [9] C. F. Manski, E. Tamer, Inference on regressions with interval data on a regressor or outcome, *Econometrica* 70 (2002) 519–546.
- 550 [10] C. Bontemps, T. Magnac, E. Maurin, Set identified linear models, *Econometrica* 80 (2012) 1129–1155.
- [11] A. Chandrasekhar, V. Chernozhukov, F. Molinari, P. Schrimpf, Inference for best linear approximations to set identified functions, CeMMAP Working Paper CWP 43/12 (2012).
- 555 [12] H. Kaido, Asymptotically efficient estimation of weighted average derivatives with an interval censored variable, *Econometric Theory* 33 (2017) 1218–1241.
- [13] K. Adusumilli, T. Otsu, Empirical likelihood for random sets, *J. Amer. Statist. Assoc.* 112 (519) (2017) 1064–1075.
- 560 [14] G. Schollmeyer, T. Augustin, Statistical modeling under partial identification: distinguishing three types of identification regions in regression analysis with interval data, *Internat. J. Approx. Reason.* 56 (part B) (2015) 224–248. doi:10.1016/j.ijar.2014.07.003.
 URL <http://dx.doi.org/10.1016/j.ijar.2014.07.003>
- 565 [15] P. Diamond, Least squares fitting of compact set-valued data, *J. Math. Anal. Appl.* 147 (1990) 351–362. doi:10.1016/0022-247X(90)90353-H.
 URL [http://dx.doi.org/10.1016/0022-247X\(90\)90353-H](http://dx.doi.org/10.1016/0022-247X(90)90353-H)
- [16] M. A. Gil, M. T. López-García, M. A. Lubiano, M. Montenegro, Regression and correlation analyses of a linear relation between random intervals, *Test* 10 (2001) 183–201.
 570

- [17] G. González-Rodríguez, Á. Blanco, N. Corral, A. Colubi, Least squares estimation of linear regression models for convex compact random sets, *Adv. Data Anal. Classif.* 1 (2007) 67–81.
- [18] B. Sinova, A. Colubi, M. Á. Gil, G. González-Rodríguez, Interval
575 arithmetic-based simple linear regression between interval data: discussion and sensitivity analysis on the choice of the metric, *Inform. Sci.* 199 (2012) 109–124. doi:10.1016/j.ins.2012.02.040.
URL <http://dx.doi.org/10.1016/j.ins.2012.02.040>
- [19] T. Maatouk, Some application of nonparametric regression with con-
580 strained data, Ph.D. thesis, University of Glasgow, Glasgow (2003).
- [20] I. Couso, D. Dubois, Statistical reasoning with set-valued information: ontic vs. epistemic views, *Internat. J. Approx. Reason.* 55 (2014) 1502–1518. doi:10.1016/j.ijar.2013.07.002.
URL <http://dx.doi.org/10.1016/j.ijar.2013.07.002>
- [21] D. F. Heitjan, D. B. Rubin, Ignorability and coarse data, *Ann. Statist.*
585 19 (4) (1991) 2244–2253.
- [22] D. F. Heitjan, Ignorability in general incomplete-data models, *Biometrika* 81 (4) (1994) 701–708.
- [23] R. D. Gill, M. J. van der Laan, J. M. Robins, Coarsening at random: Characterizations, conjectures, counter-examples, in: D. Y. Lin, T. R. Fleming (Eds.), *Proceedings of the First Seattle Symposium in Biostatistics*, Springer, New York, 1997, pp. 255–294.
590
- [24] H. T. Nguyen, On random sets and belief functions, *J. Math. Anal. Appl.* 65 (1978) 531–542.
- [25] I. Couso, D. Dubois, L. Sánchez, *Random Sets and Random Fuzzy Sets as Ill-Perceived Random Variables*, Springer, Cham, 2014.
595

- [26] M. Grabisch, Set Functions, Games and Capacities in Decision Making, Springer, Cham, 2016.
- [27] R. T. Rockafellar, Convex Analysis, Princeton University Press, Princeton,
600 1970.
- [28] R. Schneider, Convex Bodies. The Brunn–Minkowski Theory, 2nd Edition, Cambridge University Press, Cambridge, 2014.
- [29] R. A. Vitale, L_p metrics for compact, convex sets, J. Approx. Th. 45 (1985) 280–287.
- [30] H. Le, A. Kume, The Fréchet mean shape and the shape of the means,
605 Adv. Appl. Probab. 32 (2000) 101–113.
- [31] I. Molchanov, F. Molinari, Random Sets in Econometrics, Cambridge University Press, Cambridge, 2018.
- [32] E. Eisenhauer, P. Therasse, J. Bogaerts, L. Schwartz, D. Sargent, R. Ford,
610 J. Dancey, S. Arbuck, S. Gwyther, M. Mooney, L. Rubinstein, L. Shankar, L. Dodd, R. Kaplan, D. Lacombe, J. Verweij, New response evaluation criteria in solid tumours: Revised recist guideline (version 1.1), European Journal of Cancer 45 (2009) 228 – 247.
- [33] O. Gautschi, S. I. Rothschild, Q. Li, K. Matter-Walstra, A. Zippelius,
615 D. C. Betticher, M. Früh, R. A. Stahel, R. Cathomas, D. Rauch, M. Pless, S. Peters, P. Froesch, T. Zander, M. Schneider, C. Biaggi, N. Mach, A. F. Ochsenbein, Swiss Group for Clinical Cancer Research, Bevacizumab plus pemetrexed versus pemetrexed alone as maintenance therapy for patients with advanced nonsquamous non-small-cell lung cancer: Update from the
620 Swiss group for clinical cancer research (SAKK) 19/09 trial, Clin. Lung Cancer 18 (2017) 303–309.
- [34] D. R. Cox, Regression models and life-tables, J. Royal Stat. Soc., Ser. B 43 (2) (1972) 187–220.

- [35] J. Heckman, B. Singer, The identifiability of the proportional hazard model,
625 The Review of Economic Studies 51 (2) (1984) 231–241.