# Robust thin-plate splines for multivariate spatial smoothing

Ioannis Kalogridis

*Department of Mathematics, KU Leuven*

**Abstract**

A novel family of multivariate robust smoothers based on the thin-plate (Sobolev) penalty that is particularly suitable for the analysis of spatial data is proposed. The proposed family of estimators can be expediently computed even in high dimensions, is invariant with respect to rigid transformations of the coordinate axes and can be shown to possess optimal theoretical properties under mild assumptions. The competitive performance of the proposed thin-plate spline estimators relative to their non-robust counterpart is illustrated in a simulation study and a real data example involving two-dimensional geographical data on ozone concentration.

*Keywords:* Robustness, spatial data, thin-plate splines, asymptotics.
MSC 2020: 62G08, 62G35, G2H11, 62G20.

## 1. Introduction

Consider the problem of estimating the regression function $f_0 : \mathbb{R}^d \to \mathbb{R}$ from $n$ observations $(\mathbf{x}_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$, $i = 1, \ldots, n$, following the model

$$Y_i = f_0(\mathbf{x}_i) + \epsilon_i, \quad (i = 1, \ldots, n), \tag{1}$$

where the $\epsilon_i$ are random errors that are commonly assumed to be independent and identically distributed (i.i.d.) with mean zero and finite variance, but we will be able to considerably relax these assumptions over the course of this paper. Models of this general type arise naturally throughout the sciences, as very often empirical data cast doubt on parametric regression assumptions (Wood, 2017, Chapter 5).

A popular method of estimation of $f_0$ that is expounded by Wahba (1990) and Green & Silverman (1994) involves restricting $f_0$ to the multivariate Sobolev space of functions of order $m$, $\mathcal{H}^m(\mathbb{R}^d)$, i.e., the space of all functions whose partial derivatives of total order $m$ for some $m \in \mathbb{N}_+$

---

*Email address:* `ioannis.kalogridis@kuleuven.be` (Ioannis Kalogridis)

are in $\mathcal{L}^2(\mathbb{R}^d)$. Mathematically, the space $\mathcal{H}^m(\mathbb{R}^d)$ is defined as

$$\mathcal{H}^m(\mathbb{R}^d) = \left\{ f : \mathbb{R}^d \to \mathbb{R}, \ \frac{\partial^m f(x_1, \ldots, x_d)}{\partial x_1^{m_1} \ldots \partial x_d^{m_d}} \text{ exists for all } m_1 + \ldots + m_d = m \text{ and } I_m^2(f) < \infty \right\},$$

where the semi-norm $I_m : \mathcal{H}^m(\mathbb{R}^d) \to \mathbb{R}_+$ is given by

$$I_m^2(f) = \sum_{m_1+\ldots+m_d=m} \binom{m}{m_1, \ldots, m_d} \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \left( \frac{\partial^m f(x_1, \ldots, x_d)}{\partial x_1^{m_1} \ldots \partial x_d^{m_d}} \right)^2 dx_1 \ldots dx_d. \tag{2}$$

Classical least squares thin-plate spline estimators are defined as minimizers of

$$\frac{1}{2n} \sum_{i=1}^n (Y_i - f(\mathbf{x}_i))^2 + \lambda I_m^2(f),$$

over $\mathcal{H}^m(\mathbb{R}^d)$, for some $\lambda > 0$ that governs the trade-off between smoothness and goodness of fit to the data. It is easy to see that for $d = 1$ the penalty $I_m^2(f)$ reduces to $\int_{\mathbb{R}} |f^{(m)}(x)|^2 dx$ and with standard arguments this may be simplified further to $\int_{\min x_i}^{\max x_i} |f^{(m)}(x)|^2 dx$, which gives rise to classical smoothing spline estimators, see, e.g., Wahba (1990). Such simplifications are not valid for $d > 1$, but the problem still admits an elegant solution provided only that $2m > d$. In fact, for $2m > d$, the solution to this variational problem may be written in terms of $n + M$ radial and polynomial basis functions with $M = \binom{m+d-1}{d}$, see Wood (2017, p. 216) for more details. The resulting least squares thin-plate spline can be shown to converge at a fast rate under the aforementioned assumptions on the error term (Cox, 1984; Györfi et al., 2010) with the result that least squares thin-plate estimators combine computational efficiency with good theoretical properties.

An important drawback of least squares based estimators is their lack of resistance towards atypical observations. That is, these estimators tend to be overly attracted towards observations that do not follow the bulk of the data, which often leads to poor explanatory or predictive performance. In order to remedy this deficiency, this paper introduces and studies a family of generalized (M-type) thin-plate spline estimators defined as minimizers of

$$L_n(f) = \frac{1}{n} \sum_{i=1}^n \rho(Y_i - f(\mathbf{x}_i)) + \lambda I_m^2(f), \tag{3}$$

over $\mathcal{H}^m(\mathbb{R}^d)$, for some convex loss function $\rho : \mathbb{R} \to \mathbb{R}_+$. The square loss $\rho(x) = x^2/2$ is a typical example of such a loss function, but the generality of the above formulation also permits loss functions that increase less sharply as $|x| \to \infty$ thereby providing better protection against outliers. Notable examples include Huber, quantile and $L^q$ loss functions. As we discuss below, a minimizer of (3) in $\mathcal{H}^m(\mathbb{R}^d)$ exists under fairly general conditions and, as a result, to identify this

minimizer it suffices to restrict attention to an $(n + M)$-dimensional subspace spanned by radial and polynomial functions. Furthermore, we show that, for well-chosen loss functions, optimal rates of convergence may be attained without placing any moment conditions on the errors, thereby allowing for very heavy tailed error distributions within our theoretical analysis.

Thin-plate splines possess a number of notable advantages over their more popular tensor product counterparts. First, the objective function in (3) only involves one smoothing parameter irrespective of the dimension of the predictor space. By contrast, tensor product penalties require as many smoothing parameters as the number of predictors. Since smoothing parameters are usually selected from the data, tensor product smoothers entail a considerable computational burden which becomes prohibitive for higher dimensions. Another attractive property of thin-plate splines is the invariance of the penalty $I_m^2(f)$ to rotations or reflections of the coordinate axes. This fact implies that thin-plate splines are ideal for spatial/geographic smoothing as in these cases the amount of smoothing does not depend on which axes represent the relative positions, e.g., latitude and longitude. To these advantages we may add that tensor product smooths often rely on the subjective choice of the number of knots and their position, whereas with thin-plate splines the user is absolved from this responsibility, as both of these aspects emerge naturally from the mathematical problem that thin-plate splines solve.

In view of these advantages it comes as a surprise that thin-plate splines, with the exception of univariate splines with $d = 1$, have not received enough attention in the literature. In fact, both available treatments of thin-plate spline estimators of arbitrary dimensions (Cox, 1984; Györfi et al., 2010) concern the narrow class of least squares thin-plate estimators and as a result the larger and more versatile class of M-type thin-plate estimators has been overlooked. To address this gap in the literature, the present paper establishes optimal rates of convergence for a wide variety of thin-plate estimators under a mild set of conditions considerably extending our theoretical understanding. Furthermore, while the main emphasis of this work is on outlier-resistant loss functions, we also demonstrate how our methodology can be used to improve upon the results of Györfi et al. (2010) by sharpening the rate of convergence obtained by these authors for the least squares thin-plate estimator.

## 2. Main results: existence of solutions and rates of convergence

We begin our study of generalized M-type thin-plate estimators by providing sufficient conditions for the existence of a minimizer in $\mathcal{H}^m(\mathbb{R}^d)$ for the objective function $L_n(f)$, as given in (3). For the special least squares case with $\rho(x) = x^2/2$ such conditions have already been presented in the literature (see, e.g., p. 31, Wahba, 1990) and involve the uniqueness of the least squares estimator on the null-space of the penalty functional $I_m$. This null-space is $M$-dimensional with $M = \binom{m+d-1}{d}$ and consists of all polynomials of total degree at most $m$. Let $\phi_1, \ldots, \phi_M$ denote any basis for this space of polynomials. Proposition 1 below shows that the least squares requirement

for existence generalizes nicely to arbitrary convex losses.

**Proposition 1.** *Assume that $\rho$ is a convex loss function, $2m > d$ and the covariates $\mathbf{x}_i \in \mathbb{R}^d$ are such that the corresponding unpenalized M-estimator on $\phi_1, \ldots, \phi_M$ restricted to the $\mathbf{x}_i$ is unique. Then, $L_n(f)$ has a minimizer in $\mathcal{H}^m(\mathbb{R}^d)$.*

The uniqueness requirement of Proposition 1 is very mild. For example, for $d = 2$ and $m = 2$ we may take $\phi_1(x_1, x_2) = 1$, $\phi_2(x_1, x_2) = x_1$ and $\phi_3(x_1, x_2) = x_2$, so that if the estimated plane is unique then the existence of a minimizer is guaranteed. The condition $2m > d$ is satisfied for all $m \geq 1$ when $d = 1$, i.e., univariate data, but precludes small values of $m$ in higher dimensions. Unfortunately, this condition cannot be weakened as it is both necessary and sufficient for $\mathcal{H}^m(\mathbb{R}^d)$ to be a reproducing kernel Hilbert space and consequently for point evaluation maps $\mathcal{H}^m(\mathbb{R}^d) \to \mathbb{R} : f \mapsto f(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$, to be well defined. Proposition 1 is of primary importance not only as a first step towards the theoretical analysis of M-type thin-plate estimators but also for their efficient computation; Section 3 provides the details.

Having established that the problem is well-defined, we now investigate the rate of convergence of the minimizer of $L_n(f)$ in $\mathcal{H}^m(\mathcal{R}^d)$, which we denote by $\widehat{f}_n$, to the true function $f_0$. The distance metric to be used is the $\mathcal{L}^2(Q_n)$-distance given by

$$\|f - g\|_{\mathcal{L}^2(Q_n)} = \left\{ \int_{\mathbb{R}^d} |f(\mathbf{x}) - g(\mathbf{x})|^2 dQ_n(\mathbf{x}) \right\}^{1/2}$$

where $Q_n$ is the empirical measure placing mass $n^{-1}$ at each $\mathbf{x}_i \in \mathbb{R}^d$. As we shall see, under very mild conditions on the covariates, rates of convergence in terms of the empirical norm $\| \cdot \|_{\mathcal{L}^2(Q_n)}$ also translate to rates of convergence in the classical $\mathcal{L}^2$-distance. It is worth noting that, while here and in the sequel we treat the covariates as non-random, all subsequent results also hold for random covariates provided that we condition on them. Our theoretical development relies on the following regularity conditions.

A1. The loss function $\rho : \mathbb{R} \to \mathbb{R}_+$ is convex and satisfies a Lipschitz condition, i.e., there exists a $C > 0$ such that

$$|\rho(x) - \rho(y)| \leq C|x - y|, \quad \forall (x, y) \in \mathbb{R}^2.$$

A2. The errors $\epsilon_i, i = 1, \ldots, n$ are independent random variables and there exists a constant $\kappa > 0$ such that for all $|t| \leq \kappa$,

$$\inf_n \min_{1 \leq i \leq n} \mathbb{E}\{\rho(\epsilon_i + t) - \rho(\epsilon_i)\} \geq \kappa t^2$$

A3. The covariates $\mathbf{x}_i$, $i = 1, \ldots, n$ are contained in a bounded open set $\mathcal{O} \subset \mathbb{R}^d$ whose boundary,

4

$\partial\mathcal{O}$, satisfies the uniform cone condition of Adams & Fournier (2003, p. 83). That is, there exists a locally finite open cover $\{U_j\}$ of $\partial\mathcal{O}$ and a corresponding sequence of finite cones $\{C_j\}$, each congruent to some fixed cone $C$, such that

(i) There exists an $M < \infty$ such that every $U_j$ has diameter less than $M$.

(ii) $\{\mathbf{x} \in \mathcal{O} : \inf_{\mathbf{y}\in\partial\mathcal{O}} \|\mathbf{x} - \mathbf{y}\| < \delta\} \subset \bigcup_{j=1}^{\infty} U_j$ for some $\delta > 0$.

(iii) $Q_j \equiv \bigcup_{\mathbf{x}\in\mathcal{O}\cap U_j}(\mathbf{x} + C_j) \subset \mathcal{O}$ for every $j$.

(iv) For some finite $R$, every collection of $R + 1$ of the sets $Q_j$ has an empty intersection.

A4. Define the quantities

$$h_{\max,n} = \sup_{\mathbf{x}\in\mathcal{O}} \min_{1\leq i\leq n} \|\mathbf{x} - \mathbf{x}_i\|$$

$$h_{\min,n} = \min_{i\neq j} \|\mathbf{x}_i - \mathbf{x}_j\|.$$

For all large $n$, there exist finite positive constants $B_1, B_2$ such that $h_{max,n} \leq B_1$ and $h_{\max,n}/h_{\min,n} \leq B_2$.

Assumption A1 is very general and a large number of loss functions fulfils these conditions. The convexity requirement implies that $\rho$ is minimally continuous, but, importantly, differentiability is not needed for Theorem 1 below to hold. Thus, this assumption also covers non-smooth loss functions, such as the quantile loss $\rho_\alpha(x) = x(\alpha - \mathcal{I}(x < 0))$, $\alpha \in (0,1)$. Assumption A2 requires the local quadratic behaviour of $m_i(t) := \mathbb{E}\{\rho(\epsilon_i + t)\}$, $i = 1,\ldots,n$ about zero. Assumptions of this type have been extensively used in the asymptotics of M-estimators (see, e.g., van der Vaart & Wellner, 1996, Theorem 3.2.5, p. 289). It is similarly a weak condition that is satisfied quite generally. To see this, observe that by Fubini's theorem

$$m_i(t) - m_i(0) = \int_0^t \mathbb{E}\{\rho'(\epsilon_i + x)\}dx,$$

where $\rho'$ is any subgradient of $\rho$, so that all examples of loss functions given by Kalogridis (2021) also satisfy our A2 under the conditions given by that author. Our assumptions do not entail identically distributed errors, as in our experience this is too strong of an assumption for many practical settings.

Conditions A3 and A4 concerning the design points (knots) and their positions have been previously used in the asymptotics of least squares thin-plate spline estimators, see Utreras (1988). The uniform cone condition precludes sets with very irregular boundaries, but is satisfied quite generally otherwise; it is satisfied, e.g., by balls and rectangles. Condition A4 is also very modest, as it essentially requires the design points to be distinct and dense within the domain of interest, at least for large $n$. Both of these requirements follow from the fact that the ratio $h_{\max,n}/h_{\min,n}$

remains bounded for all large $n$, as if either the observations are not unique or dense enough as $n \to \infty$, $h_{\max,n}/h_{\min,n}$ could become unbounded.

With these assumptions we can now state our first asymptotic result. It is worth noting that in Theorem 1 below we treat the smoothing parameter $\lambda$ as a random variable possibly depending on our sample $(\mathbf{x}_1, Y_1), \ldots, (\mathbf{x}_n, Y_n)$. We place no restrictions on the way that $\lambda$ may depend on the sample and instead merely require that $\lambda$ is a measurable function of our sample so that $\lambda$ is a properly defined random variable. Our treatment of $\lambda$ needs to be contrasted with the treatment of the smoothing parameter by other authors, e.g., Cox (1984); Utreras (1988), who regard it as a deterministic sequence. In our view, our treatment constitutes an important extension of existing results as $\lambda$ is most often selected from the data and is thus random rather than fixed, see, e.g., the data-driven method of selecting $\lambda$ presented in Section 3.

**Theorem 1.** *Assume that the conditions of Proposition 1 are met, A1–A4 hold and further that* $\lambda = O_P(n^{-2m/(2m+d)})$ *as well as* $\lambda^{-1} = O_P(n^{2m/(2m+d)})$. *Then, there exists a sequence,* $\widehat{f}_n$, *of M-type thin-plate splines minimizing* (3) *such that*

$$\|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^2 = O_P(n^{-2m/(2m+d)}) \quad and \quad I_m(\widehat{f}_n) = O_P(1),$$

*as* $n \to \infty$.

The limit conditions involving the smoothing parameter require that $\lambda$ tends to zero in probability as $n \to \infty$, but not too fast. The rate of decay $n^{-2m/(2m+d)}$ for $\lambda$ ensures that the asymptotic variance and bias of the estimator are balanced and this leads to the rate of convergence $n^{-2m/(2m+d)}$ for thin-plate estimators. It is worth noting that $n^{-2m/(2m+d)}$ is the optimal (squared) rate of convergence for functions in $\mathcal{H}^m(\mathcal{O})$ (Stone, 1982) and therefore it cannot be improved, except in trivial cases. For $d = 1$ our rate of convergence is in agreement with the rate obtained by Kalogridis (2021) for univariate smoothing spline estimators, but for $d > 1$ the result in Theorem 1 constitutes an important generalization. The obtained rate of convergence suggests that thin-plate spline estimators, like most nonparametric estimators, suffer from the curse of dimensionality and consequently, for given sample size $n$, estimation becomes less and less precise for larger predictor dimension $d$.

Theorem 1 establishes not only a rate of convergence in the empirical norm $\| \cdot \|_{\mathcal{L}^2(Q_n)}$, but also the boundedness of the semi-norm $I_m(\widehat{f}_n)$. The latter is crucial in extending this rate of convergence to the $\mathcal{L}^2$-norm as well as in establishing optimal rates of convergence of certain useful derivatives. These extensions are given in Corollary 1 below.

**Corollary 1.** *Suppose that the assumptions of Theorem 1 hold and that* $h_{\max,n} = O(n^{-1/(2m+d)})$. *Then, the sequence of minimizers of* (3), $\widehat{f}_n$, *satisfies*

$$\int_{\mathcal{O}} |\widehat{f}_n(\mathbf{x}) - f_0(\mathbf{x})|^2 d\mathbf{x} = O_P(n^{-2m/(2m+d)}),$$

*and for every $(j_1, \ldots, j_d) \in \mathbb{R}^d$ such that $j_1 + \ldots + j_d = j \leq m$*

$$\int_{\mathcal{O}} \left| \frac{\partial \widehat{f}_n^j(\mathbf{x})}{\partial x_1^{j_1} \ldots \partial x_d^{j_d}} - \frac{\partial f_0^j(\mathbf{x})}{\partial x_1^{j_1} \ldots \partial x_d^{j_d}} \right|^2 d\mathbf{x} = O_P(n^{-2(m-j)/(2m+d)}).$$

The above rates of convergence are again optimal according to the results of Stone (1982). To the best of our knowledge, these are the first convergence results for derivatives of thin-plate estimators even in the relatively simple least squares case.

It is interesting to compare the rate of convergence $n^{-2m/(2m+d)}$ obtained herein with the rate $\log(n)n^{-2m/(2m+d)}$ for the least squares thin-plate spline estimator emerging from the results of Györfi et al. (2010, Chapter 21). The square loss $\rho(x) = x^2/2$ is not covered by our previous set of assumptions (it is not Lipschitz), but least squares estimators may be easily treated on a separate basis. In particular, letting $\widehat{f}_n$ now denote the least squares estimator, one may easily verify the inequality

$$\|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^2 + \lambda I_m^2(\widehat{f}_n) \leq 2\langle \epsilon, \widehat{f}_n - f_0 \rangle_{\mathcal{L}^2(Q_n)} + \lambda I_m^2(f_0), \tag{4}$$

where $\langle \epsilon, \widehat{f}_n - f_0 \rangle_{\mathcal{L}^2(Q_n)}$ stands for $n^{-1} \sum_{i=1}^n \epsilon_i(\widehat{f}_n(\mathbf{x}_i) - f_0(\mathbf{x}_i))$. This inequality combined with our sharper estimate for the local entropy (compare our Lemma 1 below with Lemma 20.6 in Györfi et al. (2010)) and the modulus of continuity of the empirical process derived in van de Geer (2000, Chapter 10) now leads to Theorem 2.

**Theorem 2.** *Assume that the conditions of Proposition 1 are met, A3 holds and that the errors $\epsilon_i$ are uniformly sub-Gaussian, i.e., there exist finite constants $K_1$ and $K_2 > 0$ such that*

$$\sup_n \max_{1 \leq i \leq n} K_1^2 \mathbb{E}\{e^{|\epsilon_i|^2/K_1^2} - 1\} \leq K_2.$$

*If $\lambda = O_P(n^{-2m/(2m+d)})$ and $\lambda^{-1} = O_P(n^{2m/(2m+d)})$, then there exists a sequence of least squares thin-plate spline estimators, $\widehat{f}_n$, satisfying*

$$\|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^2 = O_P(n^{-2m/(2m+d)}) \quad and \quad I_m(\widehat{f}_n) = O_P(1),$$

*as $n \to \infty$.*

It is important to emphasize that the sub-Gaussian requirement in Theorem 2 is used exclusively for the treatment of the least-squares thin-plate estimator and not for the robust estimators for the treatment of which we rely solely on assumptions A1–A4. The sub-Gaussian requirement is met in practice, e.g., whenever the errors follow a Gaussian distribution or, more generally, whenever they possess a squared exponential moment. As noted previously, Theorem 2 improves upon the corresponding result of Györfi et al. (2010). Rates of convergence for the least squares estimator in

the $\mathcal{L}^2(\mathcal{O})$-norm as well as rates of convergence for the derivatives may now be established exactly as in the proof of Corollary 1; we omit the details.

## 3. Computation and smoothing parameter selection

As hinted previously, Proposition 1 is crucial not only from a theoretical, but also from a practical standpoint. In fact, with the help of Proposition 1 and reasoning along the same lines as in the proof of Green & Silverman (1994, Theorem 7.3), it can be shown that in order to identify the minimizer of $L_n(f)$ in (3) it suffices to restrict attention to functions of the form

$$f(\mathbf{x}) = \sum_{i=1}^{n} \gamma_i \eta_{m,d}(\|\mathbf{x} - \mathbf{x}_i\|) + \sum_{j=1}^{M} \delta_j \phi_j(\mathbf{x}), \tag{5}$$

where $\eta_{m,d} : \mathbb{R}_+ \to \mathbb{R}$ is given by

$$\eta_{m,d}(x) = \begin{cases} \frac{(-1)^{m+1+d/2}}{2^{2m-1}\pi^{d/2}(m-1)!(m-d/2)!}x^{2m-d}\log(x) & d \text{ even} \\ \frac{\Gamma(d/2-m)}{2^{2m}\pi^{d/2}(m-1)!}x^{2m-d} & d \text{ odd}, \end{cases}$$

with $\Gamma(\cdot)$ denoting Euler's gamma function. The coefficient vector $\boldsymbol{\gamma}$ is subject to the set of linear constraints $\boldsymbol{\Phi}^\top \boldsymbol{\gamma} = \mathbf{0}$, where the $n \times M$ matrix $\boldsymbol{\Phi}$ has $(i,j)$th entry $\phi_j(\mathbf{x}_i)$. In other words, the coefficient vector $\boldsymbol{\gamma} \in \mathbb{R}^n$ needs to be perpendicular to the $M$-dimensional space spanned by the restriction of the polynomial functions $\phi_1, \ldots, \phi_M$ to the design points.

The representation in (5) as well as the set of linear constraints are derived from reproducing kernel Hilbert space arguments that split the space $\mathcal{H}^m(\mathbb{R}^d)$ into the closed finite-dimensional null space of $I_m(f)$ and its orthogonal complement. It follows from these arguments that plugging (5) into (3) yields the quadratic form $I_m^2(f) = \boldsymbol{\gamma}^\top \boldsymbol{\Omega} \boldsymbol{\gamma}$ where the $n \times n$ matrix $\boldsymbol{\Omega}$ has $(i,j)$th entry $\eta_{m,d}(\|\mathbf{x}_i - \mathbf{x}_j\|)$ and the minimization problem becomes

$$\min_{(\boldsymbol{\gamma}, \boldsymbol{\delta}) \in \mathbb{R}^n \times \mathbb{R}^M} \left[ \frac{1}{n} \sum_{i=1}^{n} \rho(Y_i - \boldsymbol{\omega}_i^\top \boldsymbol{\gamma} - \boldsymbol{\phi}_i^\top \boldsymbol{\delta}) + \lambda \boldsymbol{\gamma}^\top \boldsymbol{\Omega} \boldsymbol{\gamma} \right] \tag{6}$$
$$\text{s.t. } \boldsymbol{\Phi}^\top \boldsymbol{\gamma} = \mathbf{0},$$

with $\boldsymbol{\omega}_i \in \mathbb{R}^n$ and $\boldsymbol{\phi}_i \in \mathbb{R}^M$ denoting row vectors of $\boldsymbol{\Omega}$ and $\boldsymbol{\Phi}$, respectively. It can be shown that for all $\boldsymbol{\gamma} \in \mathbb{R}^n$ satisfying $\boldsymbol{\Phi}^\top \boldsymbol{\gamma} = \mathbf{0}$ we have $\boldsymbol{\gamma}^\top \boldsymbol{\Omega} \boldsymbol{\gamma} > 0$, see Wahba (1990, pp. 32–33) and the references therein. We may automatically incorporate the linear constraints and further simplify the problem by putting $\boldsymbol{\gamma} = \mathbf{Q}\boldsymbol{\beta}$ for a matrix $\mathbf{Q}$ whose columns span the null space of $\boldsymbol{\Phi}^\top$ and $\boldsymbol{\beta} \in \mathbb{R}^{n-M}$. The matrix $\mathbf{Q}$ can be easily obtained through the QR or singular value decompositions of $\boldsymbol{\Phi}$. With this simplification, the solution to (6) may be identified using, for example, the variant

of the iteratively reweighted least squares (IRLS) algorithm employed by Kalogridis (2021) with the radial basis functions and polynomials in (5) taking the place of B-spline basis functions there.

To determine the penalty parameter $\lambda$ in a data-driven way, we propose a suitable adaptation of the strategy of Maronna (2011). Let $\mathbf{r}_- = (r_{-1}, \ldots, r_{-n})^\top$ denote an approximation to the leave-one out residuals, as obtained, for example, from the last step of the IRLS algorithm. We propose to select the value of $\lambda$ that minimizes the robust cross-validation (RCV) criterion

$$\mathrm{RCV}(\lambda) = |\tau(\mathbf{r}_-)|^2,$$

where $\tau$ denotes the robust and efficient $\tau$-scale introduced by Yohai and Zamar (1988) with tuning constants equal to $c_1 = 3$ and $c_2 = 5$, corresponding to the biweighting of the mean and standard deviation respectively. This criterion may be viewed as a robustification of the celebrated leave-one-out criterion (see, e.g., Wahba, 1990, pp. 47–52) in which the $\tau$-scale is replaced by the mean of squares of the $r_{-i}$. Thus, while in the classical leave-one-out criterion all the $|r_{-i}|^2$ contribute equally (with weight $n^{-1}$ each), the use of the robust $\tau$-scale employed herein reduces the effect of large $|r_{-i}|^2$ on the selection of $\lambda$ thereby leading to a robust automatic selection procedure.

## 4. A Monte Carlo study

We now examine the practical performance of several thin-plate spline estimators by means of a simulation study. We are interested in the performance of the competing estimators not only in completely regular data, but also in data that may contain a number of outlying observations. The estimators to be considered are

- The classical least squares thin-plate spline estimator, abbreviated as LS.

- The least absolute deviations type thin-plate spline estimator with loss function $\rho(x) = |x|$ abbreviated as LAD.

- The Huber type thin-plate spline estimator with loss function

$$\rho(x) = \begin{cases} x^2/2 & |x| < 1.345 \\ 1.345|x| - 1.345^2/2 & |x| \geq 1.345 \end{cases},$$

  abbreviated as Huber.

- The logistic type thin-plate spline estimator with loss function

$$\rho(x) = 2x + 4\log(1 + e^{-x}),$$

  abbreviated as Logistic.

To compare the above estimators we generate observations from the model

$$Y_i = f_0(\mathbf{x}_i) + \epsilon_i, \quad (i = 1, \dots, n),$$

where the regression function $f_0(\mathbf{x})$ is either given by $f_1(x_1, x_2) = \exp[-8|x_1 - 0.5|^2 - 8|x_2 - 0.5|^2]$, $f_2(x_1, x_2) = \sin(2\pi x_1)\cos(2\pi x_2)$ or $f_3(x_1, x_2, x_3) = x_1^2 + x_2^2 + x_3^2$, $(x_{i1}, x_{i2}, x_{i3})$ are uniformly distributed random variables in the unit square $(0, 1)^3$ and the errors $\epsilon_i$ are generated according to the following distributions: (i) the standard Gaussian distribution, (ii) the t-distribution with 3 degrees of freedom, (iii) the skewed t-distribution with 3 degrees of freedom and non-centrality parameter equal to 1 leading to right-skewed data, (iv) a mixture Gaussian distribution with weights means equal to 1 and 10, variances equal to 1 and 0.01 and weights equal to 0.85 and 0.15 respectively and (v) Tukey's Slash distribution defined as the distribution of the quotient of independent standard normal and uniform random variables. We assess the performance of the competing estimators through the mean-squared error (MSE) given by

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} |\widehat{f}_n(\mathbf{x}_i) - f_0(\mathbf{x}_i)|^2,$$

which is an approximation to the squared $\mathcal{L}^2$-error. Table 1 below reports average MSEs and their standard errors obtained from 1000 replications with datasets of size $n = 100$.

| $f_0$ | Dist. | LS | | Huber | | Logistic | | LAD | |
|---|---|---|---|---|---|---|---|---|---|
| | | Mean | SE | Mean | SE | Mean | SE | Mean | SE |
| $f_1$ | Gaussian | 7.17 | 0.11 | 7.34 | 0.10 | 7.39 | 0.10 | 11.76 | 0.14 |
| | $t_3$ | 15.11 | 0.36 | 9.52 | 0.13 | 9.89 | 0.15 | 13.01 | 0.19 |
| | $st_{3,1}$ | 74.04 | 0.87 | 52.26 | 0.49 | 56.92 | 0.55 | 48.76 | 0.52 |
| | M. Gaussian | 292.2 | 4.22 | 21.49 | 0.40 | 36.41 | 0.68 | 21.06 | 0.39 |
| | Slash | 9e+07 | 3e+06 | 25.09 | 0.48 | 27.16 | 0.56 | 27.80 | 0.56 |
| $f_2$ | Gaussian | 14.20 | 0.16 | 14.17 | 0.17 | 13.88 | 0.17 | 21.04 | 0.23 |
| | $t_3$ | 28.25 | 0.45 | 21.98 | 0.37 | 22.02 | 0.30 | 24.44 | 0.27 |
| | $st_{3,1}$ | 84.83 | 0.90 | 62.13 | 0.55 | 66.81 | 0.58 | 59.13 | 0.57 |
| | M. Gaussian | 311.1 | 4.18 | 39.12 | 1.23 | 57.14 | 1.50 | 34.94 | 0.89 |
| | Slash | 4e+07 | 2e+06 | 72.19 | 2.63 | 56.70 | 1.38 | 57.29 | 0.95 |
| $f_3$ | Gaussian | 7.72 | 0.27 | 7.77 | 0.26 | 7.73 | 0.25 | 15.07 | 0.28 |
| | $t_3$ | 18.19 | 0.75 | 11.52 | 0.43 | 12.01 | 0.43 | 18.43 | 0.33 |
| | $st_{3,1}$ | 74.37 | 1.05 | 53.10 | 0.57 | 57.45 | 0.60 | 52.47 | 0.56 |
| | M. Gaussian | 321.90 | 6.35 | 25.89 | 0.72 | 43.02 | 1.47 | 25.09 | 0.50 |
| | Slash | 4e+07 | 5e+06 | 45.83 | 1.40 | 48.63 | 1.73 | 43.66 | 0.87 |

Table 1: Means and standard errors of 1000 MSEs ($\times 100$) with $n = 100$ of the least squares, Huber, logistic and least absolute deviations type thin-plate spline estimators.

The results in Table 1 indicate the sensitivity of least squares estimator to departures from the (sub-)Gaussian assumption. In particular, even with $t_3$ errors that lead to a few mild outliers, the least squares estimator significantly loses ground compared to the Huber, Logistic and LAD estimators. The former two estimators almost match the performance of the least squares estimator in regular data but also exhibit a high degree of resistance towards gross errors. The LAD estimator is inefficient relative to its competitors in the situation of light-tailed Gaussian errors, but exhibits superior performance to the Huber and Logistic thin-plate spline estimators in situations of heavy-tailed asymmetric contamination, such as contamination incurred by $st_{3,1}$ and mixture Gaussian errors.

To illustrate the key practical differences between the least squares thin-plate spline and its robust counterparts, Figures 1 and 2 present two examples of estimated surfaces by the least squares and least absolute deviations estimators under Gaussian and mixture Gaussian errors, respectively. The regression function $f_1$ has a unit-sized bump at $(0.5, 0.5)$, which both estimators get right in the absence of outliers, as evidenced in Figure 1. Figure 2, however, indicates that least squares estimates can be heavily distorted by the presence of outlying observations to the extent that the estimated surface bears no resemblance to the true surface. Despite the heavy contamination, the

bump is still visible in the right panel of Figure 2, which depicts the LAD estimated surface.
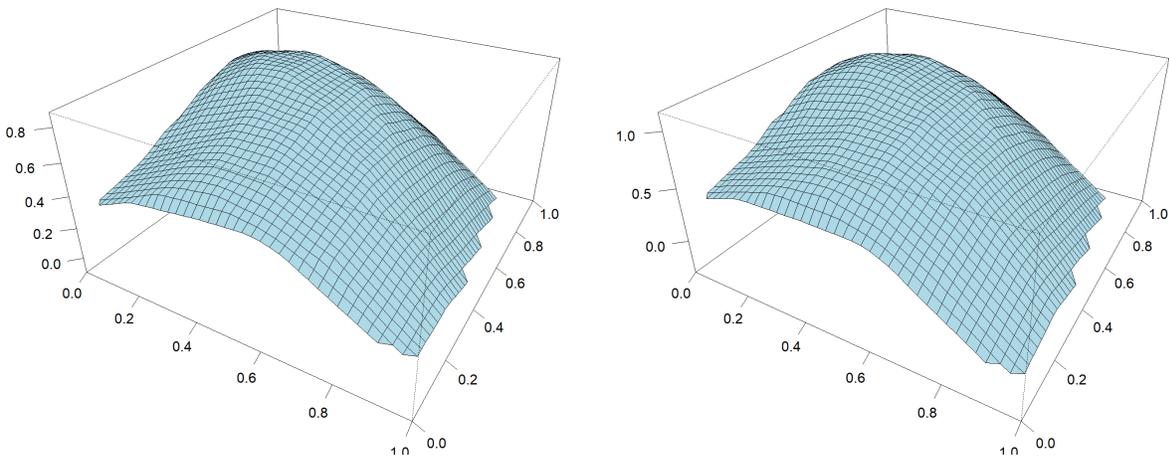


Figure 1: Typical estimated surfaces for $f_1(x_1, x_2) = \exp[-8|x_1 - 0.5|^2 - 8|x_2 - 0.5|^2]$ under Gaussian errors by the least squares and least absolute deviations estimators on the left and right panels, respectively.

## 5. Application: Ozone Levels in Midwestern USA

While stratospheric ozone protects living organisms from ultraviolet radiation from the sun, high ground-level concentrations of ozone can trigger a variety of health problems, particularly for people with breathing difficulties, children and the elderly. In this example we examine the concentration of ground-level in the midwestern US as a function of geographical longitude and latitude. That is, we consider the model

$$\text{Ozone}_i = f_0(\text{Longitude}_i, \text{Latitude}_i) + \epsilon_i, \quad (i = 1, \ldots, n).$$

The data for this analysis consists of 8-hour average surface ozone from 9am to 4pm in parts per billion (PPB) from 147 sites in midwestern US on July 3, 1987. The geographical coordinates of these sites constitute the predictor variables. The present dataset is part of a much larger dataset which is freely available as part of the `fields` package (Nychka et al., 2017) in CRAN (R Core Team, 2022).
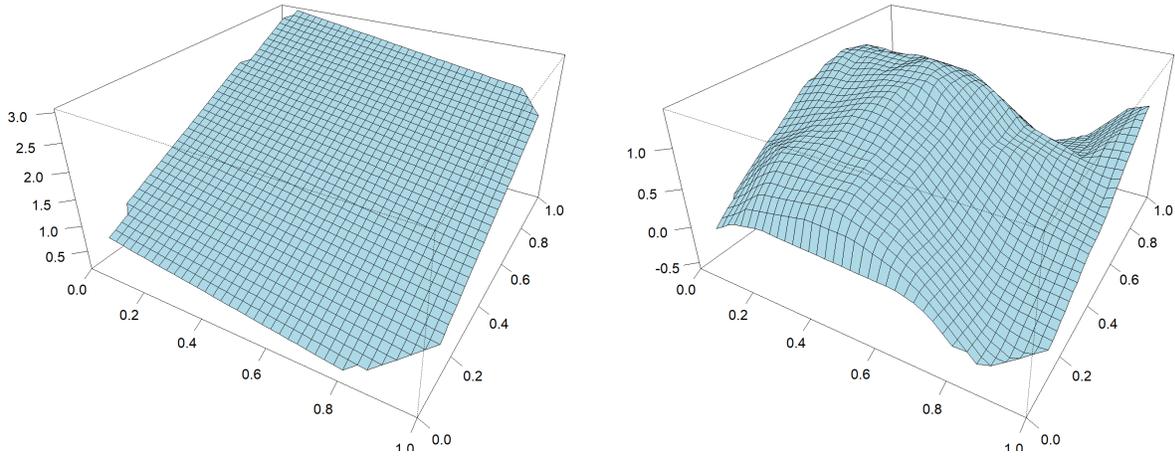
Figure 2: Typical estimated surfaces for $f_1(x_1, x_2) = \exp[-8|x_1 - 0.5|^2 - 8|x_2 - 0.5|^2]$ under mixture Gaussian errors by the least squares and least absolute deviations estimators on the left and right panels, respectively.

Since ozone concentration tends to be a skewed random variable with several potentially outlying values, we have estimated the regression function $f_0$ with both the least squares (LS) and least absolute deviations (LAD) thin-plate spline estimators. The contours of the estimated surfaces on the convex hull of the data are depicted in the left and right panels of Figure 3, respectively. The panels of the figure suggest that while there is broad agreement between these estimated surfaces around Indianapolis, the surfaces are noticeably different on the central and western part of the data. In particular, in the areas west and south of Milwaukee and St. Louis, LS estimates tend to underestimate ozone concentrations relative to LAD estimates. These differences are probably more consequential when it comes to concentrations in the excess of 60 PPB as these are more detrimental to one's health.

In view of the lack of resistance of LS estimators, it may be conjectured that the large differences between the estimates depicted in Figure 3 are attributable to the presence of atypical observations within the data. Since robust regression estimators are less attracted to outlying observations, such observations result in large absolute residuals which we may then use for their detection. A popular rule of thumb in that respect involves flagging the $i$th observation as an outlier if its standardized absolute LAD residual $r_i / \mathrm{MAD}(\mathbf{r})$ is larger than 2.5. This rule results in the detection of 14 outlying observations for the LAD estimator, but only 7 for the non-robust LS estimator.
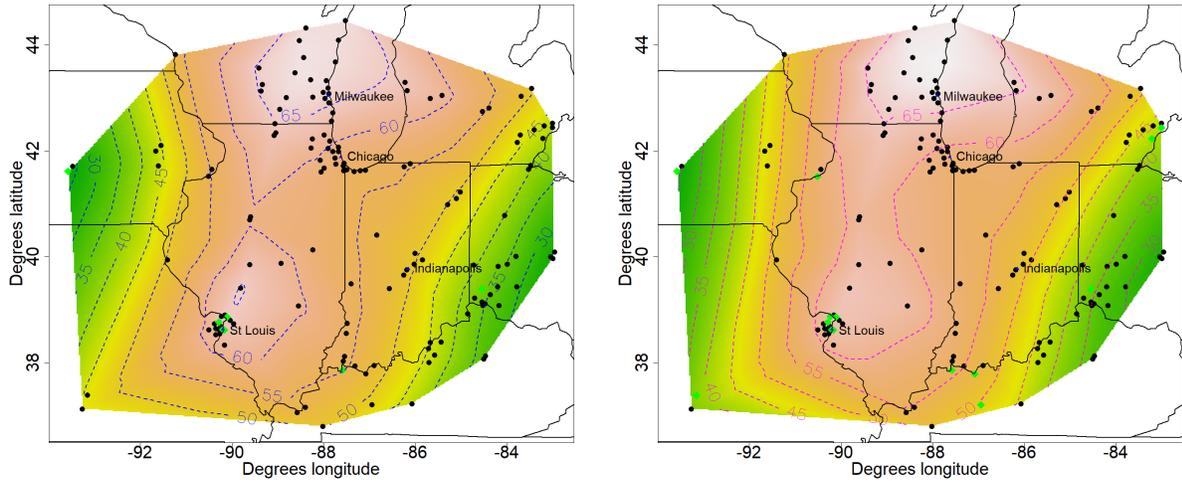
13

Figure 3: Contours of the estimated surfaces. Left: contours of the least squares thin-plate spline estimator. Right: contours of the least absolute deviations thin-plate spline estimator. Darker colors indicate higher ozone concentrations. The observation sites are depicted with solid black dots and green rhombuses depending on whether the observation is classified as an outlier or not.
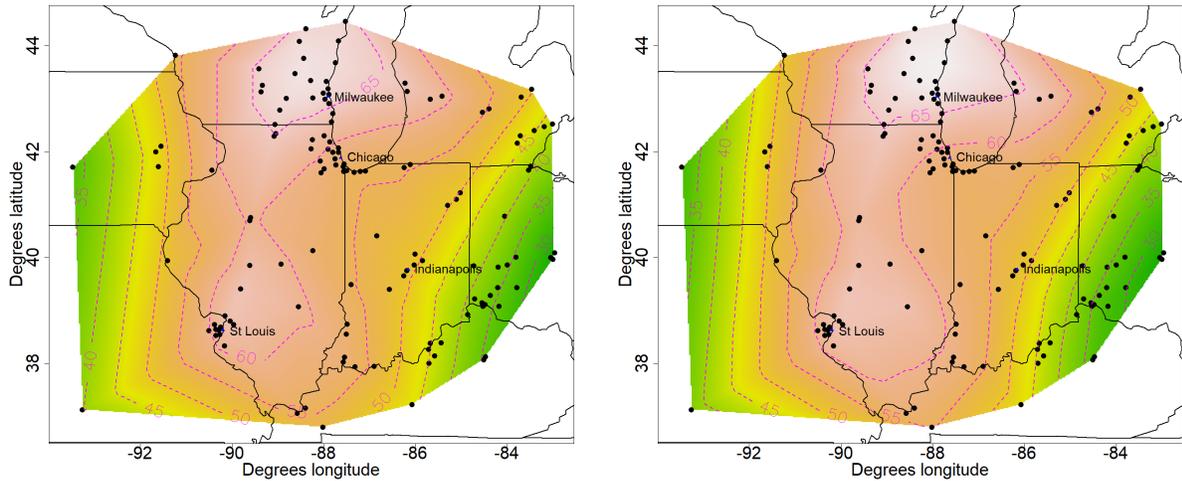


Figure 4: Contours of the estimated surfaces after removing outlying observations from the data. Left: contours of the least squares thin-plate spline estimator. Right: contours of the least absolute deviations thin-plate spline estimator. Darker colors indicate higher ozone concentrations. The observation sites are depicted with solid black dots.

To demonstrate precisely how sensitive the LS estimator is towards outlying observations, we have removed the LAD-detected outliers from the dataset and recomputed the estimates. This results in the left and right panels in Figure 4 for the LS and LAD estimates, respectively. It is interesting to observe that without these outliers the LAD and LS estimates are much closer to one

another. Moreover, comparing the right panels of Figure 3 and Figure 4 reveals that, contrary to LS estimates, LAD estimates have undergone minimal change after removing the outliers from the data. LAD estimates are thus better able to describe the bulk of the data than their sensitive LS counterparts.

## 6. Conclusion

The present paper introduces robust estimators for multivariate nonparametric regression models based on the highly advantageous thin-plate penalty. The proposed class of estimators enjoys optimal theoretical properties under mild assumptions and can be expediently computed even in high dimensions. There are several research directions worth exploring from here. Our theoretical treatment relies on a specific rate of decay of $\lambda$ and it is not known whether $\lambda$ attains this rate of decay when determined by our robust cross validation procedure. In order to ascertain whether this is the case, a close investigation of robust model selection procedures can be worthwhile. Additionally, our treatment rests upon the convexity of the loss function and hence the important class of redescending thin-plate estimators is left out. Future effort can therefore be directed towards a dedicated treatment of these estimators.

An important generalization of our ideas would involve robust estimation of multivariate non-parametric generalized linear models, where the distribution of the response variable can be any member of the exponential family, thus significantly extending the range of applications. Such an estimator may be based on, for example, the thin-plate penalty proposed herein and the density power divergence, as used by Kalogridis et al. (2023) in the univariate setting.

Another important area of research where thin-plate splines are likely to be successful is functional data analysis and, in particular, location and dispersion estimation from discretely sampled functional data. For the location case, thin-plate splines may be used to extend the robust estimator of Kalogridis and Van Aelst (2022) for discretely sampled functional data on a bounded interval to discretely sampled functional data on much more complicated multivariate domains, such as a sphere. For dispersion estimation, thin-plate splines can be combined with resistant loss functions and provide a potent alternative to estimators based on tensor product penalties, see, e.g., (Hsing & Eubank, 2015, Chapter 8). We aim to study these important generalizations as part of our future work.

## Appendix: Proofs of the theoretical results

*Proof of Proposition 1.* The result would follow by direct application of Theorem 3.2 of Cox & O'Sullivan (1985) provided that we can check the conditions of that theorem. Since $I_m^2$ is a squared semi-norm on $\mathcal{H}^m(\mathbb{R}^d)$ and $I_m$ has a finite $M$-dimensional null space, these conditions entail the

15

convexity and weak lower semicontinuity of the map

$$\mathcal{H}^m(\mathbb{R}^d) \to \mathbb{R}_+ : f \mapsto n^{-1} \sum_{i=1}^n \rho(Y_i - f(\mathbf{x}_i)).$$

We only need to check weak lower continuity, because convexity follows easily from the convexity of $\rho$. Convexity also implies that we only need to establish the lower semicontinuity of the map, as by the Hahn-Banach theorem (Rynne & Youngston, 2008) convexity and lower semicontinuity imply weak lower semicontinuity. To that end, let $\{f_k\}_k$ denote a sequence in $\mathcal{H}^m(\mathbb{R}^d)$ converging to some $f^\star$. By the Sobolev embedding theorem (Adams & Fournier, 2003, Theorem 4.12), we have

$$\max_{1 \le i \le n} |f_k(\mathbf{x}_i) - f^\star(\mathbf{x}_i)| \le c_0 \|f_k - f^\star\|_{\mathcal{H}^m(B_{r_n}^d(0))} \le c_0 \|f_k - f^\star\|_{\mathcal{H}^m(\mathbb{R}^d)},$$

for some $c_0 > 0$ not depending on $k$, where $B_{r_n}^d(0)$ denotes the smallest ball in $\mathbb{R}^d$ containing all the $\mathbf{x}_i$ and $\|\cdot\|_{\mathcal{H}^m(B_{r_n}^d(0))}, \|\cdot\|_{\mathcal{H}^m(\mathbb{R}^d)}$ denote the standard Sobolev norms on $\mathcal{B}_{r_n}^d(0)$ and $\mathbb{R}^d$, respectively. Letting $k \to \infty$, we find that $\max_{1 \le i \le n} |f_k(\mathbf{x}_i) - f^\star(\mathbf{x}_i)| \to 0$. The continuity of $\rho$ concludes the proof. $\qquad\square$

We now introduce some useful notation that will be used in the proof of Theorem 1 and Theorem 2. Let $Q_n$ denote the probability measure on $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ given by

$$Q_n(A) = \frac{1}{n} \sum_{i=1}^n \mathcal{I}(\mathbf{x}_i \in A), \ A \in \mathcal{B}(\mathbb{R}^d),$$

where $\mathcal{I}(\cdot)$ denotes the indicator function. Let $\mathcal{F}$ denote a class of real-valued functions on some domain $\mathcal{X} \subset \mathbb{R}^d$. The $\delta$-entropy for $\mathcal{F}$ in the $\mathcal{L}^2(Q_n)$-norm, $H(\delta, \mathcal{F}, \mathcal{L}^2(Q_n))$, is defined as the logarithm of the smallest number $N$ for which there exists a collection of functions $f_1, \ldots, f_N$ such that for every $f \in \mathcal{F}$ there exists a $j = j(f) \in \{1, \ldots, N\}$ with the property

$$\left\{ \int_{\mathcal{X}} |f(\mathbf{x}) - f_j(\mathbf{x})|^2 dQ_n(\mathbf{x}) \right\}^{1/2} \le \delta.$$

Similarly, we define the $\delta$-entropy with respect to the supremum norm, $H_\infty(\delta, \mathcal{F})$, as the logarithm of the smallest $N$ for which there exists $f_1, \ldots, f_N$ such that for every $f \in \mathcal{F}$ there is a $j = j(f)$ with the property

$$\sup_{\mathbf{x} \in \mathcal{X}} |f(\mathbf{x}) - f_j(\mathbf{x})| \le \delta.$$

It is clear that $H(\delta, \mathcal{F}, \mathcal{L}^2(Q_n)) \le H_\infty(\delta, \mathcal{F})$ for all $\delta > 0$. To lighten the notation in all our proofs below we will use $c_0$ to denote generic constants. Thus, the value of $c_0$ may change from appearance

to appearance.

**Lemma 1.** *Assume A3 and A4 and that $2m > d$. Then, there exists a universal constant $c_0$, independent of $n$, such that*

$$H(\delta, \{f \in \mathcal{H}^m(\mathbb{R}^d), \|f\|_{\mathcal{L}^2(Q_n)} \le 1, I_m(f) \le M\}, \mathcal{L}^2(Q_n)\} \le c_0 \left(\frac{M}{\delta}\right)^{d/m}, \quad \delta > 0, \ M \ge 1.$$

*Proof.* Observe that by A3, $Q_n$ is restricted to the open set $\mathcal{O}$ and for any $A \in \mathcal{B}(\mathbb{R}^d)$ we have $Q_n(A) = Q_n(A \cap \mathcal{O})$. Therefore, for any $(f, g) \in \mathcal{H}^m(\mathbb{R}^d) \times \mathcal{H}^m(\mathbb{R}^d)$ we find

$$\int_{\mathbb{R}^d} |f(\mathbf{x}) - g(\mathbf{x})|^2 dQ_n(\mathbf{x}) = \int_{\mathcal{O}} |f(\mathbf{x}) - g(\mathbf{x})|^2 dQ_n(\mathbf{x})$$
$$\le \sup_{\mathbf{x} \in \mathcal{O}} |f(\mathbf{x}) - g(\mathbf{x})|^2.$$

Hence, to bound $H(\delta, \mathcal{F}, \mathcal{L}^2(Q_n)\}$ it suffices to obtain a bound on

$$H_\infty(\delta, \{f \in \mathcal{H}^m(\mathcal{O}), \|f\|_{\mathcal{L}^2(Q_n)} \le 1, I_m(f) \le M\}), \quad \delta > 0, \ M \ge 1. \tag{7}$$

To bound (7) we start by showing the following inclusion

$$\{f \in \mathcal{H}^m(\mathcal{O}), \|f\|_{\mathcal{L}^2(Q_n)} \le 1, I_m(f) \le M\} \subset \{f \in \mathcal{H}^m(\mathcal{O}) : \|f\|_{\mathcal{H}^m(\mathcal{O})} \le c_0 M\}$$

for some $c_0 > 0$ and all $M \ge 1$, where $\|\cdot\|_{\mathcal{H}^m(\mathcal{O})}$ is the Sobolev norm given by

$$\|f\|_{\mathcal{H}^m(\mathcal{O})} = \left\{ \int_{\mathcal{O}} |f(\mathbf{x})|^2 d\mathbf{x} + \sum_{m_1 + \ldots + m_d = m} \binom{m}{m_1, \ldots, m_d} \int_{\mathcal{O}} \left| \frac{\partial^m f(\mathbf{x})}{\partial x_1^{m_1} \ldots \partial x_d^{m_d}} \right|^2 d\mathbf{x} \right\}^{1/2}.$$

To establish this inclusion, take an $f \in \mathcal{H}^m(\mathcal{O})$ such that $\|f\|_{\mathcal{L}^2(Q_n)} \le 1$ and $I_m(f) \le M$ for some $M \ge 1$. Our assumptions imply those of Theorem 3.4 in Utreras (1988), hence an application of that theorem reveals the existence of a constant $c_0$ such that

$$\int_{\mathcal{O}} |f(\mathbf{x})|^2 d\mathbf{x} \le c_0 \int_{\mathcal{O}} |f(\mathbf{x})|^2 dQ_n(\mathbf{x})$$
$$+ c_0 \sum_{m_1 + \ldots + m_d = m} \binom{m}{m_1, \ldots, m_d} \int_{\mathcal{O}} \left| \frac{\partial^m f(\mathbf{x})}{\partial x_1^{m_1} \ldots \partial x_d^{m_d}} \right|^2 d\mathbf{x}$$
$$\le c_0 \left\{ 1 + I_m^2(f) \right\}$$
$$\le c_0 M^2,$$

17

where the last inequality follows from our assumption that $M \geq 1$. With this bound we now obtain

$$\|f\|_{\mathcal{H}^m(\mathcal{O})} \leq \left\{ c_0 M^2 + I_m^2(f) \right\}^{1/2} \leq c_0 M,$$

as claimed.

The final step of the proof is a bound on

$$H_\infty \left( \delta, \{ f \in \mathcal{H}^m(\mathcal{O}) : \|f\|_{\mathcal{H}^m(\mathcal{O})} \leq c_0 M \} \right), \ \delta > 0, M \geq 1.$$

But this is the entropy of the closed $c_0 M$-ball and Proposition 6 of Cucker & Smale (2001) implies the existence of a universal $c_0$ such that

$$H_\infty \left( \delta, \{ f \in \mathcal{H}^m(\mathcal{O}) : \|f\|_{\mathcal{H}^m(\mathcal{O})} \leq c_0 M \} \right) \leq c_0 \left( \frac{M}{\delta} \right)^{d/m}, \ \delta > 0, M \geq 1.$$

The proof is complete.

$\square$

*Proof of Theorem 1.* The proof of the theorem employs the convexity argument of van de Geer (2002) combined with the Sobolev embedding theorem in order to localize the behaviour of the objective function around $f_0$. A tight bound on the asymptotic variance is established with the help of an exponential inequality based on the improved entropy estimates obtained in Lemma 1.

Write $L_n(f) = M_n(f) + \lambda I_m^2(f)$ where $M_n(f) = n^{-1} \sum_{i=1}^n \rho(Y_i - f(\mathbf{x}_i))$. Observe that $L_n(f)$ is the sum of two convex functions and as such it is itself convex. By Proposition 1 there exists a minimizer of $L_n(f)$ in $\mathcal{H}^m(\mathbb{R}^d)$, which we denote with $\widehat{f}_n$. Put

$$\widetilde{f}_n = \alpha \widehat{f}_n + (1 - \alpha) f_0,$$

for some $\alpha \in (0, 1)$ to be chosen. As $f_0 \in \mathcal{H}^m(\mathbb{R}^d)$ we have

$$L_n(\widetilde{f}_n) \leq \alpha L_n(\widehat{f}_n) + (1 - \alpha) L_n(f_0) \leq L_n(f_0),$$

from where, after adding $\mathbb{E}\{ M_n(\widetilde{f}_n) - M_n(f_0) \}$ on both sides, we get

$$\mathbb{E}\{ M_n(\widetilde{f}_n) - M_n(f_0) \} + I_m^2(\widetilde{f}_n) \leq \left[ M_n(f_0) - \mathbb{E}\{ M_n(f_0) \} - M_n(\widetilde{f}_n) + \mathbb{E}\{ M_n(\widetilde{f}_n) \} \right] + I_m^2(f_0). \quad (8)$$

We choose $\alpha = 1/(1 + \|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)})$. Clearly, $\alpha \in (0, 1)$ and

$$\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)} = \alpha \|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)} \leq 1.$$

Our proof consists of deriving a lower bound on the left-hand side of (8) and an upper bound on

the right-hand side of (8), both in terms of $\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}$.

We begin with the lower bound. By assumption $2m > d$ and $\mathcal{O}$ is a bounded set satisfying the uniform cone condition. Therefore, by Sobolev's theorem (Adams & Fournier, 2003, Theorem 4.12) we have the (compact) embedding

$$\mathcal{H}^m(\mathcal{O}) \to \mathcal{C}(\mathcal{O}),$$

which implies the existence of a universal $c_0 > 0$ such that for any $f \in \mathcal{H}^m(\mathcal{O}) \subset \mathcal{C}(\mathcal{O})$,

$$\sup_{\mathbf{x} \in \mathcal{O}} |f(\mathbf{x})| \leq c_0 \left\{ \int_{\mathcal{O}} |f(\mathbf{x})|^2 d\mathbf{x} + I_m^2(f) \right\}^{1/2}.$$

Approximating $\int_{\mathcal{O}} |f(\mathbf{x})|^2 d\mathbf{x}$ with $\int_{\mathcal{O}} |f(\mathbf{x})|^2 dQ_n(\mathbf{x})$, as in the proof of Lemma 1, yields

$$\sup_{\mathbf{x} \in \mathcal{O}} |f(\mathbf{x})| \leq c_0 \left\{ \int_{\mathcal{O}} |f(\mathbf{x})|^2 dQ_n(\mathbf{x}) + I_m^2(f) \right\}^{1/2}.$$

Now, since $I_m^2$ is a squared semi-norm and $I_m^2(f_0)$ is bounded, $I_m^2(f_0) \leq 1$, say, this inequality reveals that for all $f \in \mathcal{H}^m(\mathbb{R}^d)$ satisfying $\|f - f_0\|_{\mathcal{L}^2(Q_n)} \leq 1$ we have

$$\sup_{\mathbf{x} \in \mathcal{O}} |f(\mathbf{x}) - f_0(\mathbf{x})| \leq c_0 \left\{ 1 + I_m^2(f - f_0) \right\}^{1/2} \leq c_0 \{ 1 + I_m^2(f) \}^{1/2}.$$

It follows that we can choose a large enough $D_\kappa > 1$, not depending on $f$, such that

$$\sup_{\mathbf{x} \in \mathcal{O}} \frac{|f(\mathbf{x}) - f_0(\mathbf{x})|}{D_\kappa \{ 1 + I_m^2(f) \}^{1/2}} \leq \kappa,$$

where $\kappa$ is the constant in assumption A2. Therefore, for all $f \in \mathcal{H}^m(\mathbb{R}^d)$ satisfying $\|f - f_0\|_{\mathcal{L}^2(Q_n)} \leq 1$, by A2, we find

$$\begin{aligned}
\mathbb{E}\{M_n(f) - M_n(f_0)\} &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}\left\{ \rho\left(\epsilon_i + f_0(\mathbf{x}_i) - f(\mathbf{x}_i)\right) - \rho(\epsilon_i) \right\} \\
&\geq \frac{1}{n} \sum_{i=1}^n \mathbb{E}\left\{ \rho\left(\epsilon_i + \frac{f_0(\mathbf{x}_i) - f(\mathbf{x}_i)}{D_k \{ 1 + I_m^2(f) \}^{1/2}}\right) - \rho(\epsilon_i) \right\} \\
&\geq \kappa \frac{\|f - f_0\|_{\mathcal{L}^2(Q_n)}^2}{D_k^2 \{ 1 + I_m^2(f) \}}
\end{aligned}$$

19

It follows that

$$\inf_{f \in \mathcal{H}^m(\mathbb{R}^d): \|f-f_0\|_{\mathcal{L}^2(Q_n)} \leq 1} \left[ \frac{\mathbb{E}\{M_n(f) - M_n(f_0)\}}{\kappa \frac{\|f-f_0\|^2_{\mathcal{L}^2(Q_n)}}{D^2_\kappa \{1+I^2_m(f)\}}} \right] \geq 1$$

and from this, since, by construction, $\widetilde{f}_n \in \mathcal{H}^m(\mathbb{R}^d)$ and $\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)} \leq 1$, we see that

$$
\begin{aligned}
\mathbb{E}\{M_n(\widetilde{f}_n) - M_n(f_0)\} &= \frac{\mathbb{E}\{M_n(\widetilde{f}_n) - M_n(f_0)\}}{\kappa \frac{\|\widetilde{f}_n-f_0\|^2_{\mathcal{L}^2(Q_n)}}{D^2_\kappa \{1+I^2_m(\widetilde{f}_n)\}}} \kappa \frac{\|\widetilde{f}_n - f_0\|^2_{\mathcal{L}^2(Q_n)}}{D^2_\kappa \{1 + I^2_m(\widetilde{f}_n)\}} \\
&\geq \kappa \frac{\|\widetilde{f}_n - f_0\|^2_{\mathcal{L}^2(Q_n)}}{D^2_\kappa \{1 + I^2_m(\widetilde{f}_n)\}} \inf_{f \in \mathcal{H}^m(\mathbb{R}^d): \|f-f_0\|_{\mathcal{L}^2(Q_n)} \leq 1} \left[ \frac{\mathbb{E}\{M_n(f) - M_n(f_0)\}}{\kappa \frac{\|f-f_0\|^2_{\mathcal{L}^2(Q_n)}}{D^2_\kappa \{1+I^2_m(f)\}}} \right] \\
&\geq \kappa \frac{\|\widetilde{f}_n - f_0\|^2_{\mathcal{L}^2(Q_n)}}{D^2_\kappa \{1 + I^2_m(\widetilde{f}_n)\}}
\end{aligned}
\tag{9}
$$

This provides a lower bound for the left-hand side of (8) and completes the first part of our derivation.

Next, we derive an upper bound for the right-hand side of (8). To accomplish this, we need to derive the modulus of continuity of the mean-centered process $M_n(f) - \mathbb{E}\{M_n(f)\}$. We will apply Lemma 8.5 of van de Geer (2000) to this process. First, observe that by A1 for all $(f, g) \in \mathcal{H}^m(\mathbb{R}^d) \times \mathcal{H}^m(\mathbb{R}^d)$ we have

$$|\rho(Y_i - f(\mathbf{x}_i)) - \rho(Y_i - g(\mathbf{x}_i))| \leq c_0 |f(\mathbf{x}_i) - g(\mathbf{x}_i)|,$$

so that the lemma is applicable with $d_i(f, g) = |f(\mathbf{x}_i) - g(\mathbf{x}_i)|$ in the notation of van de Geer (2000). In combination with Lemma 1 we thus have

$$\sup_{f \in \mathcal{H}^m(\mathbb{R}^d): \|f-f_0\|_{\mathcal{L}^2(Q_n)} \leq 1} \left| \frac{M_n(f_0) - M_n(f) - \mathbb{E}\{M_n(f_0) - M_n(f)\}}{n^{-1/2} \|f - f_0\|^{1-d/2m}_{\mathcal{L}^2(Q_n)} \{1 + I_m(f)\}^{d/2m}} \right| = O_P(1),$$

so that

$$M_n(f_0) - M_n(\widetilde{f}_n) - \mathbb{E}\{M_n(f_0) - M_n(\widetilde{f}_n)\} = O_P(n^{-1/2}) \|\widetilde{f}_n - f_0\|^{1-d/2m}_{\mathcal{L}^2(Q_n)} \{1 + I_m(\widetilde{f}_n)\}^{d/2m}, \tag{10}$$

which provides the desired upper bound for the right-hand side of (8).

Plugging the lower bound in (9) and the upper bound in (10) into (8), we finally obtain

$$c_0 \frac{\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^2}{1 + I_m^2(\widetilde{f}_n)} + \lambda I_m^2(\widetilde{f}_n) \leq O_P(n^{-1/2})\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^{1-d/2m}\{1 + I_m(\widetilde{f}_n)\}^{d/2m} + \lambda I_m^2(f_0). \quad (11)$$

This inequality implies that

$$\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^2 = O_P(n^{-2m/(2m+d)}) \quad \text{and} \quad I_m(\widetilde{f}_n) = O_P(1). \quad (12)$$

To see this implication, note that if for real numbers $a, b, c$ we have $a \leq b + c$ then either $a \leq 2b$ or $a \leq 2c$. Applying this on (11) leads to either

$$c_0 \frac{\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^2}{1 + I_m^2(\widetilde{f}_n)} + \lambda I_m^2(\widetilde{f}_n) \leq O_P(n^{-1/2})\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^{1-d/2m}\{1 + I_m(\widetilde{f}_n)\}^{d/2m}, \quad (13)$$

or

$$c_0 \frac{\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^2}{1 + I_m^2(\widetilde{f}_n)} + \lambda I_m^2(\widetilde{f}_n) \leq 2\lambda I_m^2(f_0). \quad (14)$$

If (14) holds, (12) is easily verified. On the other hand, if (13) holds, solving it we get

$$\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)} = O_P(n^{-m/(2m+d)})\{1 + I_m(\widetilde{f}_n)\}^{\frac{d}{2m+d}}\{1 + I_m^2(\widetilde{f}_n)\}^{\frac{2m}{2m+d}},$$

as well as

$$\frac{I_m^2(\widetilde{f}_n)}{\left\{1 + I_m(\widetilde{f}_n)\right\}^{\frac{2d}{2m+d}} \left\{1 + I_m^2(\widetilde{f}_n)\right\}^{\frac{2m-d}{2m+d}}} = O_P(n^{-2m/(2m+d)})\lambda^{-1}.$$

By our assumptions on $\lambda$, the latter implies that $I_m(\widetilde{f}_n) = O_P(1)$. Hence, $\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)} = O_P(n^{-m/(2m+d)})$ verifying (12) again.

The last step in our proof involves passage from $\widetilde{f}_n$ to $\widehat{f}_n$. For this, first note that by definition of the convex combination $\widetilde{f}_n$ and (12),

$$\|\widetilde{f}_n - f_0\|_{\mathcal{L}^2(Q_n)} = \frac{\|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}}{1 + \|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}} = O_P(n^{-m/(2m+d)}),$$

whence also $\|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)} = O_P(n^{-m/(2m+d)})$. Finally, since again by (12) $I_m(\widetilde{f}_n) = O_P(1)$, by the triangle inequality we get

$$I_m(\alpha(\widehat{f}_n - f_0)) \leq I_m(\widetilde{f}_n) + I_m(f_0) = O_P(1).$$

But then also $I_m(\widehat{f}_n - f_0) = O_P(1)$, which implies the result. The proof is complete.

$\square$

*Proof of Corollary 1.* The first part of the Corollary follows from Theorem 3.4 of Utreras (1988) which in our case reads

$$\int_{\mathcal{O}} |\widehat{f}_n(\mathbf{x}) - f_0(\mathbf{x})|^2 d\mathbf{x} \le c_0 \int_{\mathcal{O}} |\widehat{f}_n(\mathbf{x}) - f_0(\mathbf{x})|^2 dQ_n(\mathbf{x})$$
$$+ c_0 h_{\max,n}^{2m} \sum_{m_1+\ldots+m_d=m} \binom{m}{m_1,\ldots,m_d} \int_{\mathcal{O}} \left| \frac{\partial^m (\widehat{f}_n(\mathbf{x}) - f_0(\mathbf{x}))}{\partial x_1^{m_1} \ldots \partial x_d^{m_d}} \right|^2 d\mathbf{x},$$

for some constant $c_0$ that does not depend on either $\widehat{f}_n$ or $f_0$. Applying now Theorem 1 and our assumption $h_{\max,n} = O(n^{-1/(2m+d)})$ yields

$$\int_{\mathcal{O}} |\widehat{f}_n(\mathbf{x}) - f_0(\mathbf{x})|^2 d\mathbf{x} = O_P(n^{-2m/(2m+d)}) + c_0 h_{\max,n}^{2m} O_P(1)$$
$$= O_P(n^{-2m/(2m+d)}) + O_P(n^{-2m/(2m+d)})$$
$$= O_P(n^{-2m/(2m+d)}),$$

as asserted.

For the second part of the Corollary we apply an interpolation inequality due to Nirenberg (1959), a simplified version of which states

$$\int_{\mathcal{O}} \left| \frac{\partial f^j(\mathbf{x})}{\partial x_1^{j_1} \ldots \partial x_d^{j_d}} \right|^2 d\mathbf{x} \le c_0 \left\{ \int_{\mathcal{O}} \left| \frac{\partial f^m(\mathbf{x})}{\partial x_1^{m_1} \ldots \partial x_d^{m_d}} \right|^2 d\mathbf{x} \right\}^{j/m} \left\{ \int_{\mathcal{O}} |f(\mathbf{x})|^2 d\mathbf{x} \right\}^{1-j/m}$$
$$+ c_0 \int_{\mathcal{O}} |f(\mathbf{x})|^2 d\mathbf{x},$$

for every $f \in \mathcal{H}^m(\mathcal{O})$, where $c_0$ is a universal constant and the inequality holds for all tuples $(j_1, \ldots, j_d)$ and $(m_1, \ldots, m_d)$ such that $j_1 + \ldots + j_d = j$ and $m_1 + \ldots + m_d = m$, respectively. Now apply this inequality with $f$ replaced by $\widehat{f}_n - f_0$ and use the first part of the corollary and Theorem 1 to get

$$\int_{\mathcal{O}} \left| \frac{\partial \widehat{f}_n^j(\mathbf{x})}{\partial x_1^{j_1} \ldots \partial x_d^{j_d}} - \frac{\partial f_0^j(\mathbf{x})}{\partial x_1^{j_1} \ldots \partial x_d^{j_d}} \right|^2 d\mathbf{x} = O_P(1) O_P(n^{-2(m-j)/(2m+d)}) + O_P(n^{-2m/(2m+d)})$$
$$= O_P(n^{-2(m-j)/(2m+d)}),$$

which completes the proof.

$\square$

*Proof of Theorem 2.* For uniformly sub-Gaussian errors $\epsilon_i$ and under Lemma 1, the derivation on

22

van de Geer (2000, p. 168) shows that

$$\frac{1}{n}\sum_{i=1}^{n}\epsilon_i(\widehat{f}_n(\mathbf{x}_i) - f_0(\mathbf{x}_i)) = O_P(n^{-1/2})\|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^{1-d/2m}\{1 + I_m(\widehat{f}_n)\}^{d/2m}.$$

Plug this into (4) to obtain

$$\|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^2 + \lambda I_m^2(\widehat{f}_n) \leq O_P(n^{-1/2})\|\widehat{f}_n - f_0\|_{\mathcal{L}^2(Q_n)}^{1-d/2m}\{1 + I_m(\widehat{f}_n)\}^{d/2m} + \lambda I_m^2(f_0).$$

Now argue as in the proof of Theorem 1 to complete the proof.

□

## Acknowledgements

## References

Adams, R.A. & Fournier, J.J.F. (2003) Sobolev Spaces, 2nd ed., Elsevier/Academic Press, Amsterdam.

Cox, D.D. (1984). Multivariate Smoothing Spline Functions. SIAM J. Numer. Anal. 21, 789–813.

Cox, D.D. & O'Sullivan, F. (1985). Analysis of penalized likelihood type estimators with applications to generalized smoothing in Sobolev Spaces. Tech. Rep. No. 51, Dept. of Statistics, University of California, Berkeley.

Cucker, F. & Smale, S. (2001). On the Mathematical Foundations of Learning. Bul. Amer. Math. Soc. 39 1–49.

Green, P.J. & Silverman, B.W. (1994). *Nonparametric Regression and Generalized Linear Models.* Chapman & Hall, London.

Györfi, L., Kohler, M., Krzyżak, A. & Walk, H. (2010). *A Distribution-Free Theory of Nonparametric Regression.* Springer, New York.

Hsing, T. & Eubank, R. (2015). *Theoretical foundations of functional data analysis, with an introduction to linear operators.* Wiley, New York.

Kalogridis, I. (2021). Asymptotics for M-type smoothing splines with non-smooth objective functions. TEST 31, 373–389.

Kalogridis, I. & Van Aelst, S. (2022). Robust optimal estimation of location from discretely sampled functional data. Scand. J. Stat., appeared online.

Kalogridis, I., Claeskens, G. & Van Aelst, S. (2023). Robust and efficient estimation of nonparametric generalized linear models. TEST, appeared online.

Maronna, R.A. (2011). Robust ridge regression for high-dimensional data. Technometrics 53, 44–53.

Nirenberg, L. (1959). On elliptic partial differential equations. Annali della Scuola Normale Superiore di Pisa - Classe di Scienze Serie 3, Tome 13, 115–162.

Nychka, D., Furrer, R., Paige, J. & Sain, S. (2017). `fields`: Tools for spatial data. `R`-package version 14.0. https://www.r-project.org/.

R Core Team (2022). R: *A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. https://www.R-project.org/.

Rynne, B. & Youngston, M.A. (2008). *Linear functional analysis*. Springer, London.

Stone, C.J. (1982). Optimal Global Rates of Convergence for Nonparametric Regression. Ann. Statist. 10, 1040–1053.

Utreras, F. (1988). Convergence Rates for Multivariate Smoothing Spline Functions. J. Approx. Theory 52, 1–27.

van de Geer, S. (2000). *Empirical processes in M-estimation*. Cambridge University Press, New York, New York.

van de Geer, S. (2002). M-estimation using penalties or sieves. J. Statist. Plann. Inference 108, 55–69.

van der Vaart, A. W. & Wellner, J. A. (1996). *Weak convergence and empirical processes*. Springer, New York.

Wahba, G. (1990). *Spline models for observational data*. Siam, Philadelphia, Pen.

Wood, S.N. (2017). Generalized Additive Models, 2nd ed.. CRC Press, Boca Raton, FL.

Yohai, V.J. & Zamar, R.H. (1988) High Breakdown-Point Estimates of Regression by Means of the Minimization of an Efficient Scale. J. Amer. Statist. Assoc. 83 406–413.