

Adaptive confidence sets in shape restricted regression

PIERRE C. BELLEC

Department of Statistics, Hill Center, Busch Campus, Rutgers University, Piscataway, NJ 08854, USA.
E-mail: pierre.bellec@rutgers.edu

A simple construction of adaptive confidence sets is proposed in isotonic, convex and unimodal regression. In univariate isotonic regression, the proposed confidence set enjoys uniform coverage over all non-decreasing regression functions. Furthermore, the diameter of the proposed confidence set automatically adapts to the unknown number of pieces of the true parameter, in the sense that the diameter is bounded from above by the minimax risk over the class of k -piecewise constant functions. The diameter of the confidence set is a simple increasing function of the number of jumps of the isotonic least-squares estimate.

A similar construction is proposed in convex regression where the true regression function is convex and piecewise affine. Here, the confidence set enjoys uniform coverage and its diameter automatically adapts to the number of affine pieces of the true regression function. The diameter of the confidence set is an increasing function of the number of affine pieces of the convex least-squares estimate.

We explain how to extend this technique to a non-convex set by proposing a similar adaptive confidence set in unimodal regression. The confidence set automatically adapts to the number of jumps of the true unimodal regression function and its diameter is an increasing function of the number of jumps of the unimodal least-squares estimate.

Keywords: adaptive confidence set; convex regression; isotonic regression; piecewise affine; piecewise constant; shape constraints; unimodal regression

1. Introduction

Let $K \subset \mathbb{R}^n$ be a closed convex set. Assume that we have the observations

$$Y_i = \mu_i + \xi_i, \quad i = 1, \dots, n,$$

where the vector $\mu = (\mu_1, \dots, \mu_n)^T \in K$ is unknown, $\xi = (\xi_1, \dots, \xi_n)^T$ is a noise vector with n -dimensional Gaussian distribution $\mathcal{N}(\mathbf{0}, \sigma^2 I_{n \times n})$ where $\sigma > 0$ and $I_{n \times n}$ is the $n \times n$ identity matrix. Denote by \mathbb{E}_μ and \mathbb{P}_μ the expectation and the probability measure corresponding to the distribution of the random variable

$$\mathbf{y} = \mu + \xi, \quad \xi \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_{n \times n}). \quad (1.1)$$

The vector $\mathbf{y} = (Y_1, \dots, Y_n)^T$ is observed and the goal is to estimate μ . Consider the scaled norm $\|\cdot\|_n$ defined by

$$\|\mathbf{u}\|_n^2 = \frac{1}{n} \sum_{i=1}^n u_i^2, \quad \mathbf{u} = (u_1, \dots, u_n)^T \in \mathbb{R}^n.$$

The error of an estimator $\hat{\mu}$ of μ is given by $\|\hat{\mu} - \mu\|_n^2$. Let $|\cdot|_2^2$ be the squared Euclidean norm, so that $\frac{1}{n}|\cdot|_2^2 = \|\cdot\|_n^2$. For a finite set E , let $|E|$ denote its cardinality. We use bold face for vectors and the components of any vector $\mathbf{v} \in \mathbb{R}^n$ are denoted by v_1, \dots, v_n . If $\mathbf{v} \in \mathbb{R}^n$ and $T \subset \{1, \dots, n\}$, we

may view \mathbf{v} as a function from $\{1, \dots, n\}$ to \mathbb{R} and denote by $\mathbf{v}_T = (v_i, i \in T) \in \mathbb{R}^T$ the restriction of this function to T . If it is clear from context, we identify the linear space \mathbb{R}^T with $\mathbb{R}^{|T|}$ with the canonical isomorphism that maps $(v_i, i \in T)$ to $(v_{k_1}, \dots, v_{k_{|T|}}) \in \mathbb{R}^{|T|}$ where $T = \{k_1, \dots, k_{|T|}\}$ with $k_1 < k_2 < \dots < k_{|T|}$.

In this paper, we consider the particular case where K is a polyhedron, that is, an intersection of a finite number of half-spaces. If the true parameter μ lies in a low-dimensional face of the polyhedron K , it has been shown that for some polyhedra K , the rate of estimation is of order $\frac{d\sigma^2}{n}$ up to logarithmic factors, where d is the dimension of the smallest face that contains μ [3,11,12,18,19,21]; see also the survey [20]. This phenomenon appears, for example, if the polyhedron K is the cone of non-decreasing sequences [3,11] or the cone of convex sequences [3,19]. For these examples, if μ lies in a d -dimensional face of the polyhedron K , the Least Squares estimator over K satisfies risk bounds and oracle inequalities with the parametric rate $\frac{d\sigma^2}{n}$, up to logarithmic factors. We consider the problem of confidence sets in this context. In particular, the present paper addresses the following questions.

- Is it possible to estimate or bound from below by a data-driven quantity the dimension d of the smallest face of the polyhedron K that contains the true parameter μ ?
- Is it possible to construct a confidence set \hat{C}_n such that:
 1. It enjoys uniform coverage over all $\mu \in K$ (i.e., $\mu \in \hat{C}_n$ with high probability).
 2. It adapts to the smallest low-dimensional face that contains μ (i.e., the diameter of \hat{C}_n should be of the order $\frac{d\sigma^2}{n}$ up to logarithmic factors if the smallest face that contains μ has dimension d).

In this paper, we answer these questions for two particular convex polyhedra, the cone of nondecreasing sequences and the cone of convex sequences, as well as for the non-convex set of unimodal sequences.

The construction of adaptive confidence sets in isotonic or convex regression has been studied in [6,7,15]. These papers show that if the true regression function is simultaneously smooth and monotone, then it is possible to construct confidence sets that adapt to the unknown smoothness of the true regression function. In the present paper, there is no smoothness assumption and the goal is to construct confidence sets that adapt to the dimension d of the smallest face of the polyhedron. A related literature studies the optimal separation rates for testing and construction of confidence sets in sparse mean models and sparse linear regression [2,8,9,26,30]. This literature is reviewed in Section 4.2 where we discuss piecewise constant signals.

The rest of the paper is organized as follows. Section 2 gives the definition of honest and adaptive confidence sets. Section 3 defines the cone of nondecreasing sequences and recalls some material from [1,24] on the statistical dimension and the intrinsic volumes of closed convex cones. In Section 4, honest and adaptive confidence sets are provided for the cone of nondecreasing sequences and Section 4.2 shows that, on the other hand, such honest and adaptive confidence sets do not exist for piecewise constant sequences. Section 4.3 discusses the case of unknown noise level. Honest and adaptive confidence sets are then given in Section 5 for the cone of convex sequences. In Section 6, we show that a similar construction is possible for the non-convex set of unimodal sequences, which shows that our construction of confidence sets does not require convexity of the model.

2. Honest and adaptive confidence sets

Let $(E_k)_{k \in J}$ be a collection of subsets K indexed by some possibly infinite set J . We will refer to the sets $(E_k)_{k \in J}$ as the *models*. Most collection of models considered in the paper will be finite and ordered by inclusion, that is,

$$E_1 \subset \dots \subset E_{k_{\max}} = K \tag{2.1}$$

where $J = \{1, \dots, k_{\text{MAX}}\}$. For any model $E_k \subset K$, the minimax risk on E_k is the quantity

$$R_{\mathbb{E}}^*(E_k) = \inf_{\hat{\mu}} \sup_{\mu \in E_k} \mathbb{E}_{\mu} \|\hat{\mu} - \mu\|_n^2,$$

where the infimum is taken over all estimators, that is, all random variables of the form $\hat{\mu} = g(\mathbf{y})$ where $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a Borel function. If $J = \{1, \dots, k_{\text{MAX}}\}$ and (2.1) holds, the minimax risks satisfy

$$R_{\mathbb{E}}^*(E_1) \leq \dots \leq R_{\mathbb{E}}^*(E_{k_{\text{MAX}}}).$$

In that case, the collection $(E_k)_{k=1, \dots, k_{\text{MAX}}}$ represents models of increasing complexity.

Similarly, if a confidence value $\alpha \in (0, 1)$ is given, one may define the minimax quantity

$$R_{\alpha}^*(E_k) = \inf \left\{ R > 0 : \sup_{\mu \in E_k} \mathbb{P}_{\mu} (\|\hat{\mu} - \mu\|_n^2 \leq R) \geq 1 - \alpha \right\} \quad (2.2)$$

for all $k \in J$, where the supremum is taken over all estimators. This quantity represents the smallest radius, in a minimax sense, of a confidence ball with confidence level $1 - \alpha$. Similarly, if $J = \{1, \dots, k_{\text{max}}\}$ and the models are ordered by inclusion as in (2.1), this quantity is an increasing function of k and we have

$$R_{\alpha}^*(E_1) \leq \dots \leq R_{\alpha}^*(E_{k_{\text{MAX}}})$$

for all $\alpha \in (0, 1)$.

The goal of this paper is to study confidence sets in shape restricted regression. A confidence set is a region \hat{C}_n such that with high probability, the unknown parameter μ belongs to \hat{C}_n . Let $\alpha \in (0, 1)$. If $\mu \in E_{k^*}$ for some $k^* \in J$, the quantity (2.2) may be used to define the oracle region

$$\hat{C}_n^*(k^*) := \{\mathbf{u} \in \mathbb{R}^n : \|\mathbf{u} - \hat{\mu}\|_n^2 \leq R_{\alpha}^*(E_{k^*})\},$$

where $\hat{\mu}$ is an estimator that achieves the supremum in (2.2) (we assume here that all infima and suprema in (2.2) are attained). Then, by definition of $R_{\alpha}^*(\cdot)$, we have that $\mu \in \hat{C}_n^*(k^*)$ with probability at least $1 - \alpha$. We call $\hat{C}_n^*(k^*)$ an *oracle* region since it is inaccessible for two reasons. First, the radius $R_{\alpha}^*(E_{k^*})$ and the integer k^* must be known in order to construct $\hat{C}_n^*(k^*)$, that is, the knowledge of the smallest model that contains μ is needed. Second, the oracle region $\hat{C}_n^*(k^*)$ is an Euclidean ball centered around the estimator $\hat{\mu}$ that achieves the infimum in (2.2), and this estimator is unknown.

This paper studies the construction of data-driven confidence sets \hat{C}_n . We consider only $1 - \alpha$ confidence sets, which means that the true parameter μ belongs to \hat{C}_n with probability at least $1 - \alpha$, uniformly over all $\mu \in K$.

We also want the diameter of the confidence set \hat{C}_n to be of the same order as the diameter of the oracle region $\hat{C}_n^*(k^*)$, that is, the value $R_{\alpha}^*(E_{k^*})$. Furthermore, the construction of \hat{C}_n should not require the knowledge of the smallest model that contains the true parameter μ : The knowledge of k^* is not needed to construct the confidence region \hat{C}_n . In that case, we say that the confidence set \hat{C}_n is adaptive.

We now give a formal definition of these properties. For any $A \subset \mathbb{R}^n$, define the diameter of A for the scaled norm $\|\cdot\|_n$ by

$$\text{diam } A := \sup_{\mathbf{v}, \mathbf{u} \in A} \|\mathbf{v} - \mathbf{u}\|_n.$$

In the following definition, for any μ we consider a probability space $(\Omega, \mathcal{A}, \mathbb{P}_{\mu})$ rich enough so that \mathbf{y} in (1.1) is a random variable. (For instance, one may take $\Omega = \mathbb{R}^n$ and \mathcal{A} the Lebesgue measure on \mathbb{R}^n .)

Definition 2.1. Let $\alpha \in (0, 1)$. Let $K \subset \mathbb{R}^n$ be a closed convex set and let $(E_k)_{k \in J}$ be a collection of subsets of K indexed by an arbitrary set J . Denote by $\hat{C}_n : \mathbb{R}^n \rightarrow \{B \subset \mathbb{R}^n\}$ a function such that (i) the event $\{\mu \in \hat{C}_n(\mathbf{y})\}$ is \mathcal{A} -measurable, and (ii) $\text{diam}(\hat{C}_n(\mathbf{y}))$ is a random variable in $(\Omega, \mathcal{A}, \mathbb{P}_\mu)$. If clear from context, we will write $\hat{C}_n = \hat{C}_n(\mathbf{y})$ for brevity. We say that \hat{C}_n is an honest confidence set if

$$\inf_{\mu \in K} \mathbb{P}_\mu(\mu \in \hat{C}_n) \geq 1 - \alpha. \quad (2.3)$$

We say that an honest confidence set \hat{C}_n is adaptive in probability if for all $\gamma \in (0, 1)$,

$$\inf_{k \in J} \inf_{\mu \in E_k} \mathbb{P}_\mu \left(\text{diam}(\hat{C}_n)^2 \leq c' R_\alpha^*(E_k) \log \left(\frac{en}{\gamma \alpha} \right)^c \right) \geq 1 - \gamma, \quad (2.4)$$

where $c' > 0$ and $c \geq 0$ are numerical constants. Alternatively to (2.4), we say that the confidence set \hat{C}_n is adaptive in expectation if for all $k \in J$,

$$\sup_{\mu \in E_k} \mathbb{E}_\mu [\text{diam}(\hat{C}_n)^2] \leq c' R_\mathbb{E}^*(E_k) \log \left(\frac{en}{\alpha} \right)^c, \quad (2.5)$$

where $c' > 0$ and $c \geq 0$ are numerical constants.

The role of the constant $c \geq 0$ is to allow for logarithmic factors. Inequality (2.3) requires that the true sequence μ lies in \hat{C}_n with probability $1 - \alpha$ for all $\mu \in K$. Inequality (2.4) implies that if the true parameter satisfies $\hat{\mu} \in E_{k^*}$ for some $k^* \in J$, then the diameter of \hat{C}_n is of the same order as the minimax quantity (2.2) of the model E_{k^*} , up to logarithmic factors.

We now consider a special case: confidence balls centered at the Least Squares estimator. The Least Squares estimator over a closed convex set K is defined by

$$\hat{\mu}^{\text{LS}}(K) = \underset{\mathbf{u} \in K}{\text{argmin}} \|\mathbf{y} - \mathbf{u}\|_n^2 = \Pi_K(\mathbf{y})$$

where Π_K denotes the convex projection onto K . By definition of the convex projection onto K , we have $(\mathbf{u} - \Pi_K(\mathbf{y}))^T(\mathbf{y} - \Pi_K(\mathbf{y})) \leq 0$ for all $\mathbf{u} \in K$, which can be rewritten as

$$\|\hat{\mu}^{\text{LS}}(K) - \mathbf{y}\|_n^2 \leq \|\mathbf{u} - \mathbf{y}\|_n^2 - \|\mathbf{u} - \hat{\mu}^{\text{LS}}(K)\|_n^2. \quad (2.6)$$

If the confidence set \hat{C}_n is an Euclidean ball, it is characterized by its center and its radius. Let $\alpha \in (0, 1)$ be a confidence value, typically $\alpha = 0.05$. Let $K \subset \mathbb{R}^n$ be a closed convex set and let $(E_k)_{k \in J}$ be a collection of subsets of K indexed by an arbitrary set J . Let \hat{r} be a positive random variable measurable with respect to \mathbf{y} and let $\hat{\mu}^{\text{LS}}(K)$ be the Least Squares estimator over K . The set

$$\hat{C}_n = \{\mathbf{v} \in \mathbb{R}^n : \|\hat{\mu}^{\text{LS}}(K) - \mathbf{v}\|_n^2 \leq \hat{r}\} \quad (2.7)$$

is an honest confidence ball if (2.3) holds. The confidence ball \hat{C}_n is said to be adaptive in probability if (2.4) holds, that is, for all $\gamma \in (0, 1)$,

$$\inf_{k \in J} \inf_{\mu \in E_k} \mathbb{P}_\mu \left(\hat{r} \leq c' R_\alpha^*(E_k) \log \left(\frac{en}{\gamma \alpha} \right)^c \right) \geq 1 - \gamma, \quad (2.8)$$

for all $\gamma \in (0, 1)$ where $c' > 0$ and $c \geq 0$ are numerical constants. The confidence ball \hat{C}_n is said to be adaptive in expectation if (2.5), that is,

$$\sup_{\mu \in E_k} \mathbb{E}_{\mu}[\hat{r}] \leq c' R_{\mathbb{E}}^*(E_k) \log\left(\frac{en}{\alpha}\right)^c, \quad (2.9)$$

for all $k \in J$, where $c' > 0$ and $c \geq 0$ are numerical constants.

3. Preliminaries

Throughout the paper, C_1, C_2, C_3, \dots denote positive absolute constants.

3.1. The cone of nondecreasing sequences and the models $(\mathcal{S}_n^{\uparrow}(k))_{k=1,\dots,n}$

Let \mathcal{S}_n^{\uparrow} be the set of all nondecreasing sequences, defined by

$$\mathcal{S}_n^{\uparrow} := \{\mathbf{u} = (u_1, \dots, u_n)^T \in \mathbb{R}^n : u_i \leq u_{i+1}, i = 1, \dots, n-1\}.$$

For $n = 1$, let $\mathcal{S}_n^{\uparrow} = \mathbb{R}$. For all $n \geq 1$, define the cone of non-increasing sequences by $\mathcal{S}_n^{\downarrow} := -\mathcal{S}_n^{\uparrow}$.

For any $\mathbf{u} \in \mathcal{S}_n^{\uparrow} \cup \mathcal{S}_n^{\downarrow}$, let $k(\mathbf{u}) := |\{u_i, i = 1, \dots, n\}|$ where $|A|$ denotes the cardinality of set A . The integer $k(\mathbf{u})$ is the smallest positive integer such that \mathbf{u} is piecewise constant with $k(\mathbf{u})$ pieces. The integer $k(\mathbf{u}) - 1$ is also the number of jumps of \mathbf{u} , that is, the number of inequalities $u_i \leq u_{i+1}$ that are strict. Define the sets

$$\mathcal{S}_n^{\uparrow}(k) = \{\mathbf{u} \in \mathcal{S}_n^{\uparrow} : k(\mathbf{u}) \leq k\}, \quad k = 1, \dots, n. \quad (3.1)$$

The set $\mathcal{S}_n^{\uparrow}(1)$ is the subspace of all constant sequences while $\mathcal{S}_n^{\uparrow}(2), \dots, \mathcal{S}_n^{\uparrow}(n-1)$ are closed non-convex sets. This defines a collection of models ordered by inclusion,

$$\mathcal{S}_n^{\uparrow}(1) \subset \mathcal{S}_n^{\uparrow}(2) \subset \dots \subset \mathcal{S}_n^{\uparrow}(n) = \mathcal{S}_n^{\uparrow}.$$

The minimax risk over the sets $\mathcal{S}_n^{\uparrow}(k)$ satisfies

$$C_4 \sigma^2 k \log \log(16n/k)/n \leq R_{\mathbb{E}}^*(\mathcal{S}_n^{\uparrow}(k)) \leq C_5 \sigma^2 k \log \log(16n/k)/n, \quad (3.2)$$

for $k \geq 2$, cf. [17]. The lower bound in (3.2) is proved in [17] using Fano's inequality, which actually yields a lower bound in probability and thus $C_6 k \log \log(16n/k)/n \leq R_{\alpha}^*(\mathcal{S}_n^{\uparrow}(k))$ for $\alpha \in (0, C_7)$. For the upper bound, Markov's inequality implies that $R_{\alpha}^*(\mathcal{S}_n^{\uparrow}(k)) \leq R_{\mathbb{E}}^*(\mathcal{S}_n^{\uparrow}(k))/\alpha$. In summary, for $R_{\alpha}^*(\mathcal{S}_n^{\uparrow}(k))$,

$$C_8(\alpha) \sigma^2 k \log \log(16n/k)/n \leq R_{\alpha}^*(\mathcal{S}_n^{\uparrow}(k)) \leq \alpha^{-1} C_9 \sigma^2 k \log \log(16n/k)/n \quad (3.3)$$

for constants that only depend on $\alpha \in (0, 1)$. Existing results also show that the dependence in α in the upper bound is at most logarithmic by properties of the Least-Squares estimator: for any $\mu \in \mathcal{S}_n^{\uparrow}(k)$, the Least-Squares estimator satisfies

$$R_{\alpha}^*(\mathcal{S}_n^{\uparrow}(k)) \leq 2\sigma^2 k \log(en/k)/n + 4\sigma^2 \log(1/\alpha)/n,$$

cf. [3], Theorem 3.2. Thus, the quantity $R_{\alpha}^*(\mathcal{S}_n^{\uparrow}(k))$ is of order $k\sigma^2/n$, up to logarithmic factors in n and $1/\alpha$.

3.2. Statistical dimension and intrinsic volumes of cones

We recall here some properties of closed convex cones. Most of the material of the present section comes from [1,24]. In the present paper, a cone is always pointed at 0. A polyhedral cone is a closed convex cone of the form

$$K = \{\mathbf{u} \in \mathbb{R}^n : \mathbf{u}^T \mathbf{v}_j \leq 0 \text{ for all } j = 1, \dots, p\}, \quad (3.4)$$

where $\mathbf{v}_1, \dots, \mathbf{v}_p$ are vectors in \mathbb{R}^n , that is, K is the intersection of a finite number of half-spaces. The polar cone of K is defined as

$$K^\circ := \{\boldsymbol{\theta} \in \mathbb{R}^n : \mathbf{v}^T \boldsymbol{\theta} \leq 0 \text{ for all } \mathbf{v} \in K\}.$$

If K is a polyhedral cone, the face of K with outward vector $\boldsymbol{\theta} \in \mathbb{R}^n$ is the set

$$F(\boldsymbol{\theta}) := \left\{ \mathbf{u} \in K : \mathbf{u}^T \boldsymbol{\theta} = \sup_{\mathbf{v} \in K} \mathbf{v}^T \boldsymbol{\theta} \right\}. \quad (3.5)$$

The face $F(\boldsymbol{\theta})$ is nonempty if and only if $\boldsymbol{\theta} \in K^\circ$. If K is the polyhedral cone (3.4) defined by the vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$, a face of a polyhedral cone K has to be of the form

$$\{\mathbf{u} \in K : \mathbf{u}^T \mathbf{v}_j = 0 \text{ for all } j \in T\}$$

for some $T \subset \{1, \dots, k\}$. The dimension of a face F is the dimension of the affine span of F .

Definition 3.1 (Statistical dimension, Amelunxen et al. [1]). For any closed convex cone $K \subset \mathbb{R}^n$, define

$$\delta(K) := \mathbb{E}[|\Pi_K(\mathbf{g})|_2^2] = \mathbb{E}[\mathbf{g}^T \Pi_K(\mathbf{g})] = \mathbb{E}\left[\left(\sup_{\boldsymbol{\theta} \in K : |\boldsymbol{\theta}|_2 \leq 1} \mathbf{g}^T \boldsymbol{\theta}\right)^2\right],$$

where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, I_{n \times n})$. The quantity $\delta(K)$ is called the statistical dimension of the cone K .

It is also known that the following holds almost surely

$$|\Pi_K(\mathbf{g})|_2^2 = \mathbf{g}^T \Pi_K(\mathbf{g}) = \left(\sup_{\boldsymbol{\theta} \in K : |\boldsymbol{\theta}|_2 \leq 1} \mathbf{g}^T \boldsymbol{\theta}\right)^2, \quad (3.6)$$

cf. [1], Proposition 3.1. The random variable (3.6) concentrates around its expectation. Indeed, the function $\mathbf{g} \mapsto |\Pi_K(\mathbf{g})|_2$ is 1-Lipschitz so that the Gaussian concentration theorem [5], Theorem 5.5, states that $|\Pi_K(\mathbf{g})|_2 \leq \mathbb{E}|\Pi_K(\mathbf{g})|_2 + \sqrt{2 \log(1/\alpha)}$ holds with probability at least $1 - \alpha$. Since $\mathbb{E}|\Pi_K(\mathbf{g})|_2 \leq \delta(K)^{1/2}$ by Jensen's inequality, this implies a concentration result for the random variable (3.6): on the same event of probability at least $1 - \alpha$, one has

$$|\Pi_K(\mathbf{g})|_2^2 \leq \delta(K) + 2\sqrt{2 \log(1/\alpha)}\delta(K) + 2 \log(1/\alpha) \leq 2\delta(K) + 4 \log(1/\alpha). \quad (3.7)$$

Our proofs also rely on the following facts given in [1], Proposition 3.1 (9), or [24], Section 5.2. If $K_1, K_2 \subset \mathbb{R}^n$ are two convex cones such that $K_1 \subset E_1$, $K_2 \subset E_2$ for two subspaces E_1, E_2 with $E_1 \perp E_2$, then $K = K_1 + K_2$ satisfies almost surely $\Pi_K(\mathbf{g}) = \Pi_{K_1}(\Pi_{E_1}(\mathbf{g}))$,

$$\Pi_K(\mathbf{g}) = \Pi_{K_1}(\mathbf{g}) + \Pi_{K_2}(\mathbf{g}), \quad |\Pi_K(\mathbf{g})|_2^2 = |\Pi_{K_1}(\mathbf{g})|_2^2 + |\Pi_{K_2}(\mathbf{g})|_2^2.$$

Consequently $\Pi_{K_1}(\mathbf{g})$ and $\Pi_{K_2}(\mathbf{g})$ are independent if $\mathbf{g} \sim N(\mathbf{0}, I_{n \times n})$ and

$$\delta(K) = \delta(K_1) + \delta(K_2). \quad (3.8)$$

If $K_1 \in \mathbb{R}^q$, $K_2 \in \mathbb{R}^{n-q}$ and $K = K_1 \times K_2$ then (3.8) still holds by applying the previous argument to the two cones $K_1 \times \{0, \dots, 0\} \subset \mathbb{R}^n$ and $\{0, \dots, 0\} \times K_2 \subset \mathbb{R}^n$.

We now define the intrinsic volumes of a polyhedral cone, which are closely related to the statistical dimension.

Definition 3.2 (Intrinsic volumes of a polyhedral cone [1]). Let $K \subset \mathbb{R}^n$ be a polyhedral cone and let $\mathbf{g} \sim N(\mathbf{0}, I_{n \times n})$. The intrinsic volumes of K are the real numbers

$$v_k(K) = \mathbb{P}(\Pi_K(\mathbf{g}) \text{ lies in the relative interior of a } k\text{-dimensional face of } K),$$

for all $k = 0, \dots, n$.

The intrinsic volumes of a polyhedral cone K define a probability distribution on the discrete set $\{0, \dots, n\}$. More precisely, define the random variable

$$V_K = \sum_{k=0}^n k \mathbf{1}_{\{\Pi_K(\mathbf{g}) \text{ lies in the relative interior of a } k\text{-dimensional face of } K\}}, \quad (3.9)$$

where $\mathbf{1}_{\{\cdot\}}$ is the indicator function. The random variable V_K is valued in $\{0, \dots, n\}$ and satisfies $\mathbb{P}(V_K = k) = v_k(K)$ for all $k = 0, \dots, n$. The following identity was derived in [1,24]:

$$\delta(K) = \sum_{k=0}^n k v_k(K), \quad (3.10)$$

that is, the statistical dimension $\delta(K)$ is the expectation of the random variable V_K . Furthermore, the random variable V_K concentrates around its expected value. The following concentration inequality is given in [24], Corollary 4.10,

$$\mathbb{P}(V_K - \delta(K) \geq \lambda) \leq \exp\left(-\frac{\delta(K)}{2} h\left(\frac{\lambda}{\delta(K)}\right)\right), \quad \text{for all } \lambda > 0,$$

where $h(t) = (1+t) \log(1+t) - t$. Using the estimate $h^{-1}(t) \leq \sqrt{2t} + 3t$ (cf. [5], Corollary 12.12), we obtain

$$\mathbb{P}(V_K - \delta(K) \geq 2\sqrt{x\delta(K)} + 6x) \leq \exp(-x), \quad \text{for all } x > 0. \quad (3.11)$$

Deriving upper and lower bounds on the statistical dimension of a cone K may be a challenging problem. Some recipes to derive such bounds are proposed in [1,10]. An exact formula is available for the statistical dimension of the cone \mathcal{S}_n^\uparrow [1], (D.12). It is given by

$$\delta(\mathcal{S}_n^\uparrow) = \delta(\mathcal{S}_n^\downarrow) = \sum_{k=1}^n \frac{1}{k}, \quad \text{so that} \quad \log n \leq \delta(\mathcal{S}_n^\uparrow) \leq \log(en). \quad (3.12)$$

Finally, we will need the following characterization of the faces of the cone \mathcal{S}_n^\uparrow .

Proposition 3.1. *Let $k \in \{1, \dots, n\}$.*

(i) *The faces of dimension k of the cone \mathcal{S}_n^\uparrow are the sets*

$$F(S) := \{\mathbf{u} = (u_1, \dots, u_n)^T \in \mathcal{S}_n^\uparrow : u_{i-1} = u_i \text{ if } i \in S\} \quad (3.13)$$

where $S \subseteq \{2, \dots, n\}$ with $|S| = n - k$. The cone \mathcal{S}_n^\uparrow has no face of dimension 0.

(ii) *For $k = 1, \dots, n$, the set $\mathcal{S}_n^\uparrow(k)$ is the union of all faces of dimension k .*

(iii) *For $K = \mathcal{S}_n^\uparrow$, the random variable V_K in (3.9) is equal to $k(\Pi_K(\mathbf{g}))$.*

Proof. (i) Let $K = \mathcal{S}_n^\uparrow$. Our definition of a face is given by an outward vector $\boldsymbol{\theta} \in \mathbb{R}^n$, and the corresponding face is given in (3.5). Every outward vector $\boldsymbol{\theta}$ gives a (possibly empty face) of K , and two distinct vectors $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2$ may produce the same face in (3.5). The face with outward vector $\boldsymbol{\theta}$ is non-empty if and only if $\sup_{\mathbf{v} \in K} \mathbf{v}^T \boldsymbol{\theta} = 0$.

Here, $\mathbf{u} \in \mathcal{S}_n^\uparrow$ if and only if $\mathbf{u} = X\mathbf{v}$ for some

$$\mathbf{v} \in \mathbb{R}^n \quad \text{with } v_2 \geq 0, \dots, v_n \geq 0 \quad (3.14)$$

and $X = (X_{ij}) \in \mathbb{R}^{n \times n}$ is the matrix with $X_{ij} = 1$ if $j \leq i$ and $X_{ij} = 0$ if $j > i$. Hence the face with outward vector $\boldsymbol{\theta}$ is non-empty if and only if $X^T \boldsymbol{\theta} \leq 0$ and $\sum_{i=1}^n \theta_i = 0$. Let $S \subset \{2, \dots, n\}$ be the set $\{i \in \{2, \dots, n\} : \theta_i + \theta_{i+1} + \dots + \theta_n < 0\}$, that is, the indices of strict inequalities in $X^T \boldsymbol{\theta} \leq 0$. Then it is clear that if $\mathbf{u} \in \mathcal{S}_n^\uparrow$ is of the form $\mathbf{u} = X\mathbf{v}$ with \mathbf{v} in (3.14), inequality $\mathbf{u}^T \boldsymbol{\theta} = 0$ implies $\mathbf{v}^T X^T \boldsymbol{\theta} = 0$ and $v_i = 0$ for $i \in S$. Since $v_i = u_i - u_{i-1}$, if the face (3.5) for outward vector $\boldsymbol{\theta}$ is non-empty, then this face is of the form (3.13). Finally, the dimension of the face (3.13) is the rank of the matrix X with rows indexed in S removed, that is, $n - |S|$.

(ii) and (iii) follow by definition of $\mathcal{S}_n^\uparrow(k)$ in (3.1) since $\Pi_K(\mathbf{g})$ belongs to the relative interior of a face $F(S)$ if and only if $F(S)$ is the smallest face of K containing $\Pi_K(\mathbf{g})$. \square

4. Nondecreasing sequences

4.1. Adaptive confidence sets for nondecreasing sequences

The estimator $\hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow)$ is the projection of \mathbf{y} onto \mathcal{S}_n^\uparrow , so the vector $\hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow)$ is nondecreasing. Let $\hat{k} = k(\hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow))$ be the number of constant pieces of the Least Squares estimator. Using this notation, we define the statistic

$$\hat{r}_\uparrow = \frac{\sigma^2}{n} \left(\sqrt{\hat{k} \log(en/\hat{k})} + \sqrt{2(\log(n/\alpha) + \hat{k} \log(en/\hat{k}))} \right)^2. \quad (4.1)$$

Theorem 4.1. *For all $\alpha \in (0, 1)$ and all $\boldsymbol{\mu} \in \mathcal{S}_n^\uparrow$, the statistic \hat{r}_\uparrow defined in (4.1) satisfies*

$$\mathbb{P}_{\boldsymbol{\mu}} \left(\left\| \hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow) - \boldsymbol{\mu} \right\|_n^2 \leq \hat{r}_\uparrow \right). \quad (4.2)$$

The above proposition shows that the confidence set (2.7) with $\hat{r} = \hat{r}_\uparrow$ satisfies condition (2.3). Up to constants and logarithmic factors, the number of constant pieces \hat{k} of the Least Squares estimator $\hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow)$ bounds the loss $\| \hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow) - \boldsymbol{\mu} \|_n^2$ from above with high probability. Since $\hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow)$ can be

computed in linear time, the integer \hat{k} and the statistic \hat{r}_\uparrow can also be computed in linear time. It is easy to compute \hat{k} visually by drawing the estimator $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow)$ and counting the number of jumps.

The proof of Theorem 4.1 relies on concentration properties of the random variable (3.6).

Proof of Theorem 4.1. Throughout the proof, we will consider partitions (T_1, \dots, T_k) of $\{1, \dots, n\}$ such that each T_j satisfies $\max T_j < \min T_{j+1}$ as well as

$$T_j = \{\min T_j, \min T_j + 1, \dots, \max T_j\}, \quad (4.3)$$

that is, T_j contains all consecutive integers from $\min T_j$ to $\max T_j$.

Let $\hat{\mu} = \hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow)$ for notational simplicity. Then (2.6) with μ replaced by $\hat{\mu}$ can be rewritten as $|\hat{\mu} - \mu|_2^2 \leq 2\xi^T(\hat{\mu} - \mu) - |\hat{\mu} - \mu|_2^2$, which implies that $|\hat{\mu} - \mu|_2 \leq \xi^T \hat{\theta}$ where $\hat{\theta} = (\hat{\mu} - \mu)/(|\hat{\mu} - \mu|_2)$ has Euclidean norm 1. By definition of $k(\cdot)$, there exists a partition $(\hat{T}_1, \dots, \hat{T}_{\hat{k}})$ of $\{1, \dots, n\}$ such that $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow)$ is constant on each \hat{T}_j , $j = 1, \dots, \hat{k}$. Since $\mu \in \mathcal{S}_n^\uparrow$, both $(\hat{\mu} - \mu)$ and $\hat{\theta}$ belong to the product cone

$$\hat{\mathcal{C}} := \mathcal{S}_{|\hat{T}_1|}^\downarrow \times \mathcal{S}_{|\hat{T}_2|}^\downarrow \times \dots \times \mathcal{S}_{|\hat{T}_{\hat{k}}|}^\downarrow. \quad (4.4)$$

By the Gaussian concentration theorem [5], Theorem 5.5, 5.6,

$$\sup_{\theta \in \mathcal{C}: |\theta|_2=1} \xi^T \theta \leq \mathbb{E} \sup_{\theta \in \mathcal{C}: |\theta|_2=1} \xi^T \theta + \sigma \sqrt{2x} \quad (4.5)$$

with probability at least $1 - e^{-x}$. Furthermore, $\mathbb{E} \sup_{\theta \in \mathcal{C}: |\theta|_2=1} \xi^T \theta \leq \sigma \delta(\mathcal{C})^{1/2}$ by Jensen's inequality, and if \mathcal{C} is of the form $\mathcal{C} = \mathcal{S}_{n_1}^\downarrow \times \dots \times \mathcal{S}_{n_k}^\downarrow$ for positive integers n_1, \dots, n_k such that $n_1 + \dots + n_k = n$, then $\delta(\mathcal{C}) = \sum_{j=1}^k \delta(\mathcal{S}_{n_j}^\downarrow) \leq \sum_{j=1}^k \log(en_j) \leq k \log(en/k)$ thanks to (3.12) and (3.8).

Let $k = 1, \dots, n$ be fixed. There are $\binom{n-1}{k}$ partitions of $\{1, \dots, n\}$ of the form (T_1, \dots, T_k) with $\max T_j < \min T_{j+1}$ for all $j = 1, \dots, k-1$ (each such partition defines a unique configuration of jumps). By the union bound and the inequality $\log \binom{n-1}{k} \leq k \log(en/k)$, we have with probability at least $1 - e^{-x}$ the bound

$$\sup_{(T_1, \dots, T_k)} \left(\sup_{\theta \in \mathcal{S}_{|T_1|}^\downarrow \times \dots \times \mathcal{S}_{|T_k|}^\downarrow: |\theta|_2=1} \xi^T \theta \right) \leq \sigma \sqrt{k \log(en/k)} + \sigma \sqrt{2(x + k \log(en/k))}.$$

Finally, we apply the union bound over all $k \in \{1, \dots, n\}$ and set $x = \log(n/\alpha)$. We have established that with probability at least $1 - \alpha$, for any random partition $(\hat{T}_1, \dots, \hat{T}_{\hat{k}})$ and $\hat{\mathcal{C}}$ in (4.4),

$$\sup_{\theta \in \hat{\mathcal{C}}: |\theta|_2=1} \xi^T \theta \leq \sigma \sqrt{\hat{k} \log(en/\hat{k})} + \sigma \sqrt{2(\log(n/\alpha) + \hat{k} \log(en/\hat{k}))}. \quad (4.6)$$

In particular, this is true for the partition $(\hat{T}_1, \dots, \hat{T}_{\hat{k}})$ induced by $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow)$ defined in the previous paragraph. Note also that the right-hand side of the previous display is equal to $(n\hat{r}_\uparrow)^{1/2}$. On this event of probability at least $1 - \alpha$ we have $|\hat{\mu} - \mu|_2 \leq \xi^T \hat{\theta} \leq (n\hat{r}_\uparrow)^{1/2}$ and the proof is complete. \square

We have established the existence of an honest confidence interval of the form

$$\hat{\mathcal{C}}_n := \{v \in \mathcal{S}_n^\uparrow : \|v - \hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow)\|_n^2 \leq \hat{r}\}.$$

This confidence set has uniform coverage over all $\mu \in \mathcal{S}_n^\uparrow$, that is, it satisfies (2.3). The next result implies that the diameter of this confidence set is minimax optimal up to logarithmic factors.

Theorem 4.2. *Let (T_1, \dots, T_k) be a partition of $\{1, \dots, n\}$ with $\max T_j < \min T_{j+1}$ for every $j = 1, \dots, k-1$, and assume that μ is constant on each T_j . Let $\gamma \in (0, 1)$. The random variable $\hat{k} = k(\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow))$ satisfies*

$$\begin{aligned}\hat{k} &\leq k \log(en/k) + 2\sqrt{\log(1/\gamma)k \log(en/k)} + 6\log(1/\gamma) \\ &\leq 2k \log(en/k) + 7\log(1/\gamma)\end{aligned}\quad (4.7)$$

with probability greater than $1 - \gamma$ as well as

$$\mathbb{E}_\mu[\hat{k}] \leq k \log(en/k). \quad (4.8)$$

Furthermore $\sup_{\mu \in \mathcal{S}_n^\uparrow(k)} \mathbb{E}_\mu[\hat{k}] \geq k \log(n/k)$ when k divides n so that (4.8) is sharp.

Note that in Theorem 4.2 μ is only assumed to be piecewise constant, but not necessarily monotone. This will be useful in Section 6.

Proof of Theorem 4.2. Define the closed convex cone

$$K = \mathcal{S}_{|T_1|}^\uparrow \times \mathcal{S}_{|T_2|}^\uparrow \times \dots \times \mathcal{S}_{|T_k|}^\uparrow \subset \mathbb{R}^n \quad (4.9)$$

and let $\hat{\mu}^* = \Pi_K(\mathbf{y})$. It is clear that

$$\min_{\mathbf{u} \in K} \sum_{j=1}^k \|\mathbf{y}_{T_j} - \mathbf{u}_{T_j}\|_2^2 = \min_{\mathbf{u}_1 \in \mathcal{S}_{|T_1|}^\uparrow, \dots, \mathbf{u}_k \in \mathcal{S}_{|T_k|}^\uparrow} \sum_{j=1}^k \|\mathbf{y}_{T_j} - \mathbf{u}_{T_j}\|_2^2 = \sum_{j=1}^k \min_{\mathbf{u}_j \in \mathcal{S}_{|T_j|}^\uparrow} \|\mathbf{y}_{T_j} - \mathbf{u}_{T_j}\|_2^2.$$

Thus, as $\mathbf{y} = \mu + \xi$ and μ is constant on each T_j , we have

$$\hat{\mu}_{T_j}^* = \Pi_{\mathcal{S}_{|T_j|}^\uparrow}(\mathbf{y}_{T_j}) = \Pi_{\mathcal{S}_{|T_j|}^\uparrow}(\mu_{T_j} + \xi_{T_j}) = \mu_{T_j} + \Pi_{\mathcal{S}_{|T_j|}^\uparrow}(\xi_{T_j}).$$

As adding the constant sequence μ_{T_j} does not modify the number of constant pieces (or the number of jumps), we have

$$k(\hat{\mu}_{T_j}^*) = k(\mu_{T_j} + \Pi_{\mathcal{S}_{|T_j|}^\uparrow}(\xi_{T_j})) = k(\Pi_{\mathcal{S}_{|T_j|}^\uparrow}(\xi_{T_j})) = k((\Pi_K(\xi))_{T_j}).$$

Let V_K be the random variable defined in (3.9). By the properties of product cones given in [24], Section 5.2, V_K has the same distribution as

$$V_{\mathcal{S}_{|T_1|}^\uparrow} + \dots + V_{\mathcal{S}_{|T_k|}^\uparrow}$$

and the random variables $V_{\mathcal{S}_{|T_1|}^\uparrow}, \dots, V_{\mathcal{S}_{|T_k|}^\uparrow}$ are independent (cf. [24], (5.2)). By Proposition 3.1, for all $j = 1, \dots, k$ we have $k(\hat{\mu}_{T_j}^*) = V_{\mathcal{S}_{|T_j|}^\uparrow}$ so that V_K and $\sum_{j=1}^k k(\hat{\mu}_{T_j}^*)$ have the same distribution. By (3.10), $\mathbb{E}V_K = \delta(K)$ and by (3.11), with probability greater than $1 - \gamma$ we have

$$V_K \leq \delta(K) + 2\sqrt{\log(1/\gamma)\delta(K)} + 6\log(1/\gamma).$$

To bound $\delta(K)$ from above, we use that the statistical dimension of a direct product of cones is the sum of the statistical dimensions (cf. [1], Proposition 3.1, or the discussion preceding (3.8)),

$$\delta(K) = \sum_{j=1}^k \delta(\mathcal{S}_{|T_j|}^\uparrow) \leq \sum_{j=1}^k \log(e|T_j|) \leq k \log(en/k),$$

where we have used (3.12) and Jensen's inequality for the two inequalities.

The random variable V_K is distributed as $\sum_{j=1}^k k(\hat{\mu}_{T_j}^*)$. Thus, to complete the proof, it is enough to prove that almost surely, $\hat{k} := k(\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow)) \leq \sum_{j=1}^k k(\hat{\mu}_{T_j}^*)$. Let $\hat{\mu} = \hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow)$ for notational simplicity. It is clear that

$$k(\hat{\mu}) = |\{\hat{\mu}_i, i = 1, \dots, n\}| \leq \sum_{j=1}^k k(\hat{\mu}_{T_j}) = \sum_{j=1}^k |\{\hat{\mu}_i, i \in T_j\}|,$$

since a piece counted on the left-hand side must be counted at least once (and possibly multiple times) on the right-hand side. For all $j = 1, \dots, k$, the vectors $\hat{\mu}_{T_j}$ and $\hat{\mu}_{T_j}^*$ are solutions of the minimization problems

$$\hat{\mu}_{T_j}^* = \underset{\mathbf{v} \in \mathcal{S}_{|T_j|}^\uparrow}{\operatorname{argmin}} \|\mathbf{v} - \mathbf{y}_{T_j}\|_2^2, \quad \hat{\mu}_{T_j} = \underset{\substack{\mathbf{v} \in \mathcal{S}_{|T_j|}^\uparrow: \\ \hat{\mu}_{\min(T_j)} \leq \mathbf{v}_1, \\ \mathbf{v}_{|T_j|} \leq \hat{\mu}_{\max(T_j)}}}{\operatorname{argmin}} \|\mathbf{v} - \mathbf{y}_{T_j}\|_2^2.$$

This means that $\hat{\mu}_{T_j}$ solves a minimization problem with additional constraints at the boundary. By Lemma B.1, we have

$$k(\hat{\mu}_{T_j}) \leq k(\hat{\mu}_{T_j}^*)$$

for all $j = 1, \dots, k$, which completes the proof of Equation (4.8).

To show that $\sup_{\mu \in \mathcal{S}_n^\uparrow(k)} \mathbb{E}_\mu[k] \geq k \log(n/k)$ when k divides n , consider a partition (T_1, \dots, T_k) with $|T_j| = n/k$ and $\max T_j \leq \min T_{j+1}$ and set $\mu_i = xj$ if $i \in T_j$ for some parameter $x > 0$. The jumps of μ have amplitude x . In the event $x > 2 \max_{i=1}^n |\xi_i|$ we have $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow) = \hat{\mu}^{\text{LS}}(K)$ where K is the cone (4.9). By dominated convergence and using the fact that adding a constant vector to $\mathbf{u} \in \mathbb{R}^{|T_j|}$ does not change the number of jumps of $\Pi_{\mathcal{S}_{|T_j|}^\uparrow}(\mathbf{u})$,

$$\lim_{x \rightarrow +\infty} \mathbb{E}_\mu[k(\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow))] = \mathbb{E}_\mu[k(\hat{\mu}^{\text{LS}}(K))] = \mathbb{E}_\mu[k(\hat{\mu}^{\text{LS}}(K))] = \delta(K) = \sum_{j=1}^k \delta(\mathcal{S}_n^\uparrow(|T_j|)).$$

This quantity is greater than $k \log(n/k)$ due to the lower bound in (3.12). \square

Corollary 4.3. *Let $J = \{1, \dots, n\}$ and define the collection of models $(E_k)_{k \in J} = (\mathcal{S}_n^\uparrow(k))_{k \in J}$. The random variable \hat{r}_\uparrow defined in (4.1) satisfies (4.2), (2.8) and (2.9) with \hat{r} replaced by \hat{r}_\uparrow . Thus, the ball centered at $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow)$ of radius $\sqrt{\hat{r}_\uparrow}$ is an honest confidence set, which is adaptive in probability and in expectation with respect to the models $(\mathcal{S}_n^\uparrow(k))_{k=1, \dots, n}$.*

Proof. To show that (2.9) holds, for any $\mu \in \mathcal{S}_n^\uparrow(k)$ we have

$$\mathbb{E}_\mu[\hat{r}_\uparrow]^{1/2} \leq \sigma n^{-1/2} (\sqrt{2 \log(n/\alpha)} + (1 + \sqrt{2}) \sqrt{\mathbb{E}_\mu[\hat{k} \log(en/\hat{k})]})$$

thanks to $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ and the triangle inequality for $\mathbb{E}_{\mu}[(\cdot)^2]^{1/2}$. Since $x \mapsto x \log(en/x)$ is concave on $[1, n]$, $\mathbb{E}_{\mu}[\hat{k} \log(en/\hat{k})] \leq \mathbb{E}_{\mu}[\hat{k}] \log(en/\mathbb{E}_{\mu}[\hat{k}])$. The bound (4.8) then yields $\mathbb{E}_{\mu}[\hat{k}] \times \log(en/\mathbb{E}_{\mu}[\hat{k}]) \leq k \log(en/k)^2$ which completes the proof of (2.9) for the minimax risk given in (3.2), since (2.9) allows logarithmic factors.

To show that (2.8) holds for the minimax risk in (3.3), we instead use (4.7). This yields that with probability at least $1 - \gamma$,

$$\begin{aligned}\hat{r}_\uparrow^{1/2} &\leq \sigma n^{1/2} (\sqrt{2 \log(n/\alpha)} + (1 + \sqrt{2}) \sqrt{\hat{k} \log(en/\hat{k})}) \\ &\leq \sigma n^{1/2} (\sqrt{2 \log(n/\alpha)} + (1 + \sqrt{2}) \sqrt{2k \log(en/k) \log(en) + 7 \log(1/\gamma) \log(en)}).\end{aligned}$$

This completes the proof since (2.8) allows logarithmic factors. \square

4.2. Adaptive confidence sets do not exist for piecewise constant signal

The existence of honest and adaptive confidence sets with respect to the models $(\mathcal{S}_n^\uparrow(k))_{k=1,\dots,n}$ is surprising given the negative results for seemingly similar models. For the problem of constructing adaptive confidence sets in sparse linear regression with p covariates and $n \leq p$ observations, where one observes

$$Y_i = \sum_{j=1}^p X_{ij} \beta_j^* + \xi_i, \quad i = 1, \dots, n \text{ where } X_{ij}, \xi_i \sim \text{iid } \mathcal{N}(0, 1),$$

Nickl and van de Geer [26], Theorem 2, shows the impossibility to construct confidence sets for β^* that are (a) honest over all sparsity levels $k_1 \in \{k = 1, 2, 3, \dots : k \leq n/\log p\}$ and (b) adaptive over sparsity levels $k_0 \in \{k = 1, 2, 3, \dots : k_1 \leq \sqrt{n}/\log p\}$. This impossibility result holds as $n, p, k \rightarrow +\infty$ with $|\beta^*|_0 = p^{1-\eta}$ for some constant $\eta \in (0, 1)$. Due to this impossibility result, confidence sets over all sparsity levels $k_1 \in \{k = 1, 2, 3, \dots : k \leq n/\log p\}$ cannot have squared radius smaller than the minimax estimation rate $|\beta^*|_0 \log(p)/n$ for $|\beta^*|_0 \leq \sqrt{n}/\log p$. This impossibility result for honest and adaptive confidence sets in sparse linear regression builds upon impossibility results for signal detection [22]. Similar impossibility results have been obtained for testing sparsity levels in the Gaussian mean model [2,9] and for testing sparsity levels in linear regression [8]. We refer the reader to Carpenter and Verzelen [9] and Table 1 therein for a complete characterization of the different phase transitions for testing sparsity levels in the Gaussian mean model, and [8], Table 1, for the corresponding phase transitions in sparse linear regression.

The models $(\mathcal{S}_n^\uparrow(k))_{k=1,\dots,n}$ can be related to sparse linear regression in the following way:

$$\mathcal{S}_n^\uparrow(k) = \{X\mathbf{v}, \mathbf{v} \in \mathbb{R} \times \mathbb{R}_+^{n-1} \text{ with } \|\mathbf{v}\|_0 = k - 1\} \quad (4.10)$$

and $X \in \mathbb{R}^{n \times n}$ is the design matrix with entries $X_{ij} = I_{i \leq j}$. Here, $\mathbb{R}_+ = [0, +\infty)$ and $I_{i \leq j}$ is the indicator function equal to 1 if and only if $i \leq j$. Thus, the model $\mathcal{S}_n^\uparrow(k)$ can be seen as a sparse linear regression model with design X and at most k nonzero entries, with an additional non-negativity constraint.

The major difference in (4.10) compared to sparse linear regression is the non-negativity of v_2, \dots, v_n in (4.10), which encodes the shape constraint. If one drops the non-negativity constraints, the above model becomes the set of piecewise constant signals with $k - 1$ jumps:

$$\mathcal{B}_n(k) = \{\mathbf{v} \in \mathbb{R}^n : |\{i = 2, \dots, n : v_i \neq v_{i-1}\}| = k - 1\}.$$

The minimax rate over $\mathcal{B}_n(k)$ satisfies

$$R_{\mathbb{E}}^*(\mathcal{B}_n(k)) \asymp \begin{cases} \sigma^2/n & \text{if } k = 1, \\ k \log \log(16n/k) \sigma^2/n & \text{if } k = 2, \\ k \log(en/k) \sigma^2/n & \text{if } k \geq 3, \end{cases} \quad (4.11)$$

cf. [17], Theorem 4.1. The following result shows that honest and adaptive confidence sets do not exist for the models $(\mathcal{B}_n(k))_{k=1,\dots,n}$. This a consequence of the lower bound in [2].

Proposition 4.4. *Let $\alpha \in (0, 1/8)$ independent of n . In the model (1.1), let \hat{C}_n be a confidence set such that, for all $k \leq n^{1/2}$ inequality $\inf_{\mu \in \mathcal{B}_n(k)} \mathbb{P}_{\mu}(\mu \in \hat{C}_n) \geq 1 - \alpha$ holds. Then there exists an absolute constant C_{12} such that*

$$\mathbb{E}_{\mathbf{0}}[\text{diam}(\hat{C}_n)^2] \geq \sigma^2 C_{12} n^{-1/2} \quad (4.12)$$

for all n larger than some absolute constant $C_{10} > 0$.

Thus a confidence set \hat{C}_n that is honest with respect to the models $(\mathcal{B}_n(k))_{k \leq \sqrt{n}}$ must incur a squared radius of order at least $\sigma^2 n^{-1/2}$ which is larger than the minimax rate in (4.11) for $k \ll \sqrt{n}$. Furthermore, the discrepancy between the best possible squared radius and the minimax rate (4.11) is polynomial in n when $k = \lfloor n^{\kappa} \rfloor$ for $\kappa < 1/2$:

$$\frac{\mathbb{E}_{\mathbf{0}}[\text{diam}(\hat{C}_n)^2]}{R_{\mathbb{E}}^*(\mathcal{B}_n(\lfloor n^{\kappa} \rfloor))} \geq C_{11}(\kappa) \frac{n^{1/2-\kappa}}{\log n}$$

for all confidence sets that are honest over $(\mathcal{B}_n(k))_{k \leq \sqrt{n}}$. The constant in the previous display depends on κ only. The $\sigma^2 n^{-1/2}$ lower bound in (4.12) also appears for confidence sets in sparse regression [8, 22, 26] and testing sparsity levels in sparse mean models [2, 9]. The non-decreasing shape constraint of the previous section leads to a fundamentally different behavior: Theorems 4.1 and 4.2 show that it is possible to construct confidence sets that are both honest and adaptive (up to logarithmic factors) over all models in (3.1), i.e., with squared radius of the same order as the minimax rate.

Because any $k/2$ -sparse vector is also piecewise constant with at most k jumps, Proposition 4.4 is a straightforward consequence of the lower bound given in [2], Theorem 1.

Proof of Proposition 4.4. Let $r_n^2 = 8\mathbb{E}_{\mathbf{0}}[\text{diam}(\hat{C}_n)^2]$. We follow the argument laid out in [26], Section 3.1, to relate confidence sets to signal detection problems. Let $k_1 = \lfloor \sqrt{n} \rfloor$ be the largest integer smaller or equal to \sqrt{n} and consider testing

$$H_0 : \mu = \mathbf{0}, \quad \text{against } H_1 : \mu \in \mathcal{B}_n(k_1) \text{ with } \|\mu\|_n^2 \geq r_n^2.$$

With a minor abuse of notation, we denote by H_1 the set $H_1 = \{\mu \in \mathcal{B}_n(k_1) \text{ with } \|\mu\|_n^2 \geq r_n^2\}$. Consider the test statistic T_n valued in $\{0, 1\}$ defined by $T_n = 1$ if $\hat{C}_n \cap H_1 \neq \emptyset$ and $T_n = 0$ otherwise. The type I and type II errors of T_n are bounded as follows:

- Under the null, on the event $\Omega_0 = \{\mathbf{0} \in \hat{C}_n\} \cap \{\text{diam}(\hat{C}_n)^2 < 8\mathbb{E}_{\mathbf{0}}[\text{diam}(\hat{C}_n)^2]\}$ we have $\hat{C}_n \cap H_1 = \emptyset$ so that $\mathbb{P}_{\mathbf{0}}(T_n = 1) \leq \mathbb{P}(\Omega_0^c) \leq \alpha + 1/8$ by the union bound and Markov's inequality.
- Under the alternative, if $\mu \in H_1$ then $\mathbb{P}_{\mu}(T_n = 0) \leq \mathbb{P}_{\mu}(\mu \notin \hat{C}_n) \leq \alpha$.

The sum of type I and type II errors of the above test is thus at most $2\alpha + 1/8$. Define now $H'_1 = \{\boldsymbol{\mu} \in \mathbb{R}^n : |\boldsymbol{\mu}|_0 \leq k_1/2\}$. Since any $(k_1/2)$ -sparse vector is also piecewise constant with k_1 pieces, $H'_1 \subset H_1$ holds and

$$\begin{aligned} \mathbb{P}_0(T_n = 1) + \sup_{\boldsymbol{\mu} \in H'_1} \mathbb{P}_{\boldsymbol{\mu}}(T_n = 0) &\leq \mathbb{P}_0(T_n = 1) + \sup_{\boldsymbol{\mu} \in H_1} \mathbb{P}_{\boldsymbol{\mu}}(T_n = 0) \\ &\leq 2\alpha + 1/8 \leq 3/8. \end{aligned}$$

By [2], Theorem 1, if $r_n^2 \leq C_{12}\sigma^2 n^{-1/2}$ for some absolute constant $C_{12} > 0$, there exists no test T_n with sum of type I and type II errors less than $3/8$ for testing $H_0 : \boldsymbol{\mu} = 0$ against H'_1 . This implies (4.12) and completes the proof. \square

4.3. Adaptivity to unknown noise level in non-decreasing sequences

Although the radius \hat{r}_\uparrow in (4.1) depends on σ^2 , Theorem 4.1 is also helpful to construct estimators $\hat{\sigma}^2$ of σ^2 . Indeed, Theorem 4.1 implies that with probability at least $1 - \alpha - \beta$,

$$\begin{aligned} |1 - \sigma^{-1} \|\mathbf{y} - \hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow)\|_n| &\leq |1 - \sigma^{-1} \|\boldsymbol{\xi}\|_n| + \sigma^{-1} \|\hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow) - \boldsymbol{\mu}\|_n \\ &\leq \sqrt{2/n} (1 + \sqrt{\log(1/\beta)}) + \sigma^{-1} \hat{r}_\uparrow^{1/2} \end{aligned} \quad (4.13)$$

where we used that $\boldsymbol{\xi} \rightarrow \|\boldsymbol{\xi}\|_n - \sigma$ is a $n^{-1/2}$ -Lipschitz function of $\boldsymbol{\xi}$ with $\mathbb{E}[\|\boldsymbol{\xi}\|_n - \sigma]^2 \leq 2\sigma^2(1 - \sqrt{1 - 1/n}) \leq 2\sigma^2/n$ and the Gaussian concentration result [5], Theorem 5.5, which gives $\|\boldsymbol{\xi}\|_n - \sigma \leq \sigma\sqrt{2/n} + \sigma\sqrt{2\log(1/\beta)/n}$. Now let

$$\hat{\epsilon} = n^{-1/2} (\sqrt{2} + \sqrt{2\log(1/\beta)} + \sqrt{\hat{k}\log(en/\hat{k})} + \sqrt{2(\log(n/\alpha) + \hat{k}\log(en/\hat{k}))})$$

be the right-hand side in (4.13), which only depends on α, \hat{k}, n . Set also $\hat{\sigma}^2 = \|\mathbf{y} - \hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow)\|_n^2$. Then on the above event of probability at least $1 - 2\alpha$, we have

$$|1 - \hat{\sigma}/\sigma| \leq \hat{\epsilon}, \quad \sigma(1 - \hat{\epsilon})_+ \leq \hat{\sigma} \leq \sigma(1 + \hat{\epsilon}), \quad (1 + \hat{\epsilon})^{-1} \hat{\sigma} \leq \sigma \leq (1 - \hat{\epsilon})_+^{-1} \hat{\sigma}$$

so that for instance $(1 - \hat{\epsilon})^{-1} \hat{\sigma}$ provides a fully data-driven, high-probability upper bound for σ as soon as $\hat{\epsilon}$ is bounded away from 1. We obtain, as a consequence, the following data-driven analog of the results of the previous section with fully data-driven radius.

Theorem 4.5 (Unknown σ). *Let $\alpha, \beta \in (0, 1)$ be fixed constants independent of n . Consider the above notation for $\hat{\epsilon}$ and define*

$$\tilde{r}_\uparrow = \frac{\|\mathbf{y} - \hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow)\|_n^2}{(1 - \hat{\epsilon})_+^2 n} (\sqrt{\hat{k}\log(en/\hat{k})} + \sqrt{2(\log(n/\alpha) + \hat{k}\log(en/\hat{k}))})^2.$$

Then we have

- (i) $\mathbb{P}(\|\hat{\boldsymbol{\mu}}^{\text{LS}}(\mathcal{S}_n^\uparrow) - \boldsymbol{\mu}\|^2 \leq \tilde{r}_\uparrow) \geq 1 - \alpha - \beta$.
- (ii) $k(\boldsymbol{\mu})/n \rightarrow 0$ as $n \rightarrow +\infty$ implies $\hat{\epsilon} \rightarrow 0$ in probability.
- (iii) $k(\boldsymbol{\mu})/n \rightarrow 0$ as $n \rightarrow +\infty$ implies $\tilde{r}_\uparrow = O_{\mathbb{P}}(\sigma^2(k(\boldsymbol{\mu})/n) \log(en/k(\boldsymbol{\mu})))^2$.

Hence the fully data-driven radius \tilde{r}_\uparrow is honest and shrinks adaptively with respect to the models $(\mathcal{S}_n^\uparrow(k))_{k \leq na_n}$ for any sequence $a_n \rightarrow 0$. Alternatively, one may use the median absolute deviation to estimate σ^2 as explained, for instance, in Remark 3.1 of [17].

Proof of Theorem 4.5. (i) follows from Theorem 4.1 and (4.13). For (ii), if $a_n = k(\mu)/n$ converges to 0 then by Theorem 4.2 with $\gamma = 1/n$ we have $\hat{k}/n \leq 2a_n \log(en^{-1}) + C_{13} \log(n)/n =: b_n$ with probability at least $1/n$. Furthermore, $a_n \rightarrow 0$ implies $b_n \rightarrow 0$. On the same event we have $\hat{\epsilon} \leq (1 + \sqrt{2})\{b_n \log(en^{-1})\}^{1/2} + C_{14}\sqrt{\log(en)/n}$ which also converges to 0. For (iii), $\hat{\epsilon} \rightarrow 0$ in probability implies that $\hat{\sigma}^2(1 - \hat{\epsilon})_+^{-2} = O_{\mathbb{P}}(\sigma^2)$, so that and Theorem 4.2 provide the required upper bound on \tilde{r}_\uparrow . \square

5. Adaptive confidence sets for convex sequences

Confidence sets can also be obtained in univariate convex regression. If $n \geq 3$, define the set of convex sequences \mathcal{S}_n^\cup by

$$\mathcal{S}_n^\cup := \{\mathbf{u} = (u_1, \dots, u_n)^T \in \mathbb{R}^n : 2u_i \leq u_{i+1} + u_{i-1}, i = 2, \dots, n-1\},$$

and define $\mathcal{S}_n^\cup = \mathbb{R}$ if $n = 1$ and $\mathcal{S}_n^\cup = \mathbb{R}^2$ if $n = 2$. For all $n \geq 1$, define the cone of concave sequences by $\mathcal{S}_n^\cap := -\mathcal{S}_n^\cup$.

For any $\mathbf{u} \in \mathcal{S}_n^\cup$, let $q(\mathbf{u}) - 1 \geq 0$ be the number of inequalities $2u_i \leq u_{i+1} + u_{i-1}$, $i = 2, \dots, n-1$ that are strict. The integer $q(\mathbf{u})$ is also the smallest positive integer such that \mathbf{u} is piecewise affine with $q(\mathbf{u})$ pieces. Define the sets

$$\mathcal{S}_n^\cup(q) = \{\mathbf{u} \in \mathcal{S}_n^\cup : q(\mathbf{u}) \leq q\}, \quad q = 1, \dots, n-1.$$

The set $\mathcal{S}_n^\cup(1)$ is the subspace of all affine sequences while $\mathcal{S}_n^\cup(2), \dots, \mathcal{S}_n^\cup(n-2)$ are closed non-convex sets. We have

$$\mathcal{S}_n^\cup(1) \subset \mathcal{S}_n^\cup(2) \subset \dots \subset \mathcal{S}_n^\cup(n-1) = \mathcal{S}_n^\cup.$$

These sets represent models of increasing complexity.

There exist numerical constants $c, c' > 0$ such that for all $\alpha \leq (0, \min(c, 1))$ and any $q = 1, \dots, n-1$, we have

$$\frac{c'\sigma^2 q}{n} \leq R_\alpha^*(\mathcal{S}_n^\cup(q)) \leq \frac{16\sigma^2 q \log(en/q)}{n} + \frac{4\log(1/\alpha)}{n}, \quad (5.1)$$

cf. [3], Theorem 4.3, for the upper bound and [4], Proposition 7, for the lower bound. Thus, for $\alpha > 0$ small enough, the quantity $R_\alpha^*(\mathcal{S}_n^\cup(q))$ is of order $q\sigma^2/n$, up to logarithmic factors.

The statistical dimension of the cone \mathcal{S}_n^\cup satisfies [3], Theorem 4.1,

$$\delta(\mathcal{S}_n^\cup) = \delta(\mathcal{S}_n^\cap) \leq 8\log(en). \quad (5.2)$$

It is not known whether this upper bound is sharp. However, the fact that the statistical dimension of \mathcal{S}_n^\cup grows slower than a logarithmic function of n is enough for the purpose of the present paper.

The following bound on the risk of $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\cup)$ will be useful.

Proposition 5.1. *Let $\mu \in \mathcal{S}_n^{\cup}$. Then*

$$\mathbb{E}_{\mu} \|\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\cup}) - \mu\|_2^2 \leq \mathbb{E}_{\mu} \left[\left(\sup_{\mathbf{v} \in \mathcal{T}_{\mu}: \|\mathbf{v}\|_2 \leq 1} \xi^T \mathbf{v} \right)^2 \right] = \sigma^2 \delta(\mathcal{T}_{\mu}) \leq 8\sigma^2 q(\mu) \log \frac{en}{q(\mu)},$$

where \mathcal{T}_{μ} is the tangent cone at μ defined by $\mathcal{T}_{\mu} := \{\mathbf{u} - t\mu, \mathbf{u} \in \mathcal{S}_n^{\cup}, t \in \mathbb{R}\}$.

An outline of the proof of this result is as follows. More details may be found in [3], Proposition 4.3.

Outline of the proof of Proposition 5.1. The inequality $\mathbb{E}_{\mu} \|\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\cup}) - \mu\|_2^2 \leq \sigma^2 \delta(\mathcal{T}_{\mu})$ was proved by [27] and it is a direct consequence of (2.6) with $\mathbf{u} = \mu$. To bound from above the statistical dimension of \mathcal{T}_{μ} , we have the inclusion

$$\mathcal{T}_{\mu} \subset \mathcal{S}_{|T_1|}^{\cup} \times \cdots \times \mathcal{S}_{|T_q(\mu)|}^{\cup},$$

where $(T_1, \dots, T_{q(\mu)})$ is a partition of $\{1, \dots, n\}$ such that μ is affine on each T_j , $j = 1, \dots, q(\mu)$. The formula for the statistical dimension of a direct product of cones (cf. [1], Proposition 3.1, or the discussion preceding (3.8)) yields

$$\delta(\mathcal{S}_{|T_1|}^{\cup} \times \cdots \times \mathcal{S}_{|T_q(\mu)|}^{\cup}) = \sum_{j=1}^{q(\mu)} \delta(\mathcal{S}_{|T_j|}^{\cup}) \leq 8 \sum_{j=1}^{q(\mu)} \log(e|T_j|) \leq 8 \log(en/q(\mu)),$$

where we used (5.2) and Jensen's inequality. \square

We now turn to the construction of confidence sets. Recall that if $\mathbf{u} \in \mathcal{S}_n^{\cup}$ is a convex sequence, $q(\mathbf{u})$ is the number of pieces in the piecewise affine decomposition of \mathbf{u} . Let $\hat{q} := q(\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\cup}))$ be the number of affine pieces of the Least Squares estimator. Then, define the statistic

$$\hat{r}_{\cup} = \frac{\sigma^2 \hat{q} (16 + 24 \log(n) + 4 \log(1/\alpha))}{n}. \quad (5.3)$$

Similarly to the case of the statistic \hat{r}_{\uparrow} in isotonic regression, the following result shows that the confidence ball (2.7) with $\hat{r} = \hat{r}_{\cup}$ enjoys uniform coverage over all $\mu \in \mathcal{S}_n^{\cup}$.

Theorem 5.2. *For all $\alpha \in (0, 1)$ and all $\mu \in \mathcal{S}_n^{\cup}$, the statistic \hat{r}_{\cup} defined in (4.1) satisfies*

$$\|\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\cup}) - \mu\|_n^2 \leq \hat{r}_{\cup}, \quad (5.4)$$

with probability at least $1 - \alpha$.

The above result is analogous to Theorem 4.1. The numerical constants are slightly worse in the case of the present section because the upper bound (5.2) on the statistical dimension of the cone \mathcal{S}_n^{\cup} is slightly worse than (3.12). The proof of Theorem 5.2 is similar to the proof of Theorem 4.1 and can be found in Appendix C.

Now, the goal is to show that the statistic \hat{r}_{\cup} is of the same order as the minimax quantity (5.1). We employ a different strategy than in the previous section. For any function $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ which is weakly differentiable, the divergence of g is the random variable

$$D_g(\mathbf{y}) = \sigma^2 \sum_{i=1}^n \frac{\partial}{\partial y_i} g(\mathbf{y})_i.$$

It is well known that by Stein's identity, under suitable conditions on g (cf. [25], Section 2, or [29], Lemma 3.6), we have

$$\sigma^2 \mathbb{E}_\mu D_g(\mathbf{y}) = \mathbb{E}_\mu [\xi^T g(\mathbf{y})]. \quad (5.5)$$

The divergence of the estimator $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\cup) = \Pi_{\mathcal{S}_n^\cup}(\mathbf{y})$ is given in [14], Proposition 2.7, (see also [25]). Namely, we have the following result.

Proposition 5.3 ([14,25]). *If $g(\cdot) = \Pi_{\mathcal{S}_n^\cup}(\cdot)$ is the projection onto the cone of convex sequences, then (5.5) holds and we have*

$$D_g(\mathbf{y}) = \hat{q} + 1$$

almost surely, where $\hat{q} = q(\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\cup))$.

This result can be used to bound from above the expected radius of the statistic \hat{r}_\cup .

Theorem 5.4. *Let $\mu \in \mathcal{S}_n^\cup$. Then*

$$\mathbb{E}_\mu [\hat{q}] \leq 8q(\mu) \log \frac{en}{q(\mu)} - 1. \quad (5.6)$$

Consequently, for all $\alpha \in (0, 1)$, the statistic (5.3) satisfies

$$\mathbb{E}_\mu [\hat{r}_\cup] \leq \frac{\sigma^2 q(\mu) \text{polylog}(n, 1/\alpha)}{n}. \quad (5.7)$$

where $\text{polylog}(n, 1/\alpha) = 8 \log(en)(16 + 24 \log(n) + 4 \log(1/\alpha))$.

Proof. By Proposition 5.3 and (5.5), we have

$$\sigma^2 \mathbb{E}_\mu [1 + \hat{q}] = \mathbb{E}_\mu [\xi^T \Pi_{\mathcal{S}_n^\cup}(\mathbf{y})] = \mathbb{E}_\mu [\xi^T (\Pi_{\mathcal{S}_n^\cup}(\mathbf{y}) - \mu)]. \quad (5.8)$$

By the Cauchy–Schwarz inequality, we have

$$\begin{aligned} \sigma^2 \mathbb{E}_\mu [1 + \hat{q}] &\leq \mathbb{E}_\mu^{1/2} \left[\left(\frac{\xi^T (\Pi_{\mathcal{S}_n^\cup}(\mathbf{y}) - \mu)}{|\Pi_{\mathcal{S}_n^\cup}(\mathbf{y}) - \mu|_2} \right)^2 \right] \mathbb{E}_\mu^{1/2} |\Pi_{\mathcal{S}_n^\cup}(\mathbf{y}) - \mu|_2^2, \\ &\leq \sigma \sqrt{\delta(\mathcal{T}_\mu)} \mathbb{E}_\mu^{1/2} |\Pi_{\mathcal{S}_n^\cup}(\mathbf{y}) - \mu|_2^2. \end{aligned}$$

Using Proposition 5.1 completes the proof of (5.6). Inequality (5.7) is a direct consequence of (5.6) and of the definition of \hat{r}_\cup . \square

The above result is different from Theorem 4.2 in isotonic regression. Theorem 4.2 controls both the expectation and the deviations of \hat{k} . In this section, Theorem 5.4 only controls the expectation of \hat{q} . This comes from the use of Stein's identity in the proof of Theorem 5.4, which yields a result only in expectation.

The arguments used to prove Theorem 4.2 are based on the concentration properties of the intrinsic volumes of cones, while the proof of Theorem 5.4 relies on Stein's identity and Proposition 5.3. Thus, we have presented two methods to bound from above the expected diameter of the confidence sets constructed in the present paper. The concentration inequalities used in Theorem 4.2 for non-decreasing sequences provides exponential deviation inequalities, which are much stronger than the

bounds in expectation obtained in Theorem 5.4 (Indeed, bounds in expectation only imply weak deviation bounds by Markov's inequality). On the other hand, the argument based on Stein's formula in Theorem 5.4 has wider applicability. For instance, it readily applies to \hat{k} from Section 4: Since $\mathbb{E}_{\mu}[\hat{k}]\sigma^2 = \mathbb{E}_{\mu}[\xi^T(\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\uparrow}) - \mu)]$ holds [25], we have

$$\mathbb{E}_{\mu}[\hat{k}] \leq \mathbb{E}[|\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\uparrow}) - \mu|_2^2]^{1/2} \sigma \delta(\mathcal{T}_{\mu})^{1/2} \leq \sigma^2 k(\mu) \log(en/k(\mu))$$

for $\mu \in \mathcal{S}_n^{\uparrow}$ by the same argument as in (5.8) for the first inequality and [3], Theorem 3.2, for the second, $\mathcal{T}_{\mu} = \{\mathbf{u} - t\mu, (t, \mathbf{u}) \in \mathbb{R} \times \mathcal{S}_n^{\uparrow}\}$ is the tangent cone. This provides an alternative proof of (4.8).

We end this section with the following consequence of Theorems 5.2 and 5.4.

Corollary 5.5. *Let $J = \{1, \dots, n-1\}$ and define the collection of models $(E_k)_{k \in J} = (\mathcal{S}_n^{\cup}(k))_{k=1, \dots, n-1}$. The random variable \hat{r}_{\cup} defined in (5.3) satisfies (5.4) and (2.9) with \hat{r} replaced by \hat{r}_{\cup} . Thus, the ball centered at $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\cup})$ of radius $\sqrt{\hat{r}_{\cup}}$ is an honest confidence set, which is adaptive in expectation with respect to the models $(\mathcal{S}_n^{\cup}(k))_{k=1, \dots, n-1}$.*

The proof follows by combining Theorems 5.2 and 5.4, similarly to the proof of Corollary 4.3 for non-decreasing sequences.

6. Non-convexity: Adaptive confidence sets for unimodal sequences

Let $m \in \{1, \dots, n\}$. A sequence $\mathbf{u} \in \mathbb{R}^n$ is unimodal with mode at position m if and only if $\mathbf{u}_{\{1, \dots, m\}}$ is non-increasing and $\mathbf{u}_{\{m, \dots, n\}}$ is nondecreasing—in other words, \mathbf{u} belongs to the set

$$K_m := \{\mathbf{u} = (u_1, \dots, u_n)^T \in \mathbb{R}^n : u_1 \geq \dots \geq u_m \leq u_{m-1} \leq \dots \leq u_n\}.$$

Next, by taking the union over all possible locations for the mode, we define the set of all unimodal sequences as

$$\mathcal{U} := \bigcup_{m=1, \dots, n} K_m.$$

The set \mathcal{U} is non-convex for $n \geq 3$.

Recall that if $\mathbf{u} \in \mathcal{S}_n^{\uparrow} \cup \mathcal{S}_n^{\downarrow}$ is monotone, $k(\mathbf{u}) = |\{u_i, i = 1, \dots, n\}|$ is the smallest positive integer k such that \mathbf{u} is piecewise constant constant with k pieces. We extend the function k to the set of unimodal sequences by setting

$$\forall \mathbf{u} \in \mathcal{U}, \quad k(\mathbf{u}) = \min\{k \geq 1 : \exists \text{ partition } (T_1, \dots, T_k) : \mathbf{u}_{T_j} \text{ is constant for all } j = 1, \dots, k\},$$

where the partition (T_1, \dots, T_k) is a partition of $\{1, \dots, n\}$ with $\max T_j < \min T_{j+1}$ for all $j = 1, \dots, k-1$. A unimodal sequence $\mathbf{u} \in \mathcal{U}$ has $k(\mathbf{u}) - 1$ jumps. Similarly to isotonic regression, define the models

$$\mathcal{U}(k) = \{\mathbf{u} \in \mathcal{U} : k(\mathbf{u}) \leq k\}, \quad k = 1, \dots, n. \quad (6.1)$$

These sets define models of increasing complexity since

$$\mathcal{U}(1) \subset \mathcal{U}(2) \subset \dots \subset \mathcal{U}(n) = \mathcal{U}$$

It is known that there exist numerical constants c, c', c'' such that for all $\alpha \leq c$,

$$c'\sigma^2 k/n \leq R_\alpha^*(\mathcal{U}(k)) \leq c''\sigma^2(k \log(en/k) + \log(n/\alpha))/n. \quad (6.2)$$

Indeed, the lower bound is a consequence of (3.3) and the inclusion $R_\alpha^*(\mathcal{S}_n^\uparrow(k)) \subset R_\alpha^*(\mathcal{U}(k))$, while the upper bound is proved in [3,13,16], see [3], Appendix C, for explicit constants. Thus, for $\alpha > 0$ small enough, the quantity $R_\alpha^*(\mathcal{S}_n^\uparrow(k))$ is of order $k\sigma^2/n$, up to logarithmic factors in n and $1/\alpha$. Similarly, the minimax risk over the sets $\mathcal{U}(k)$ satisfies

$$c'\sigma^2 k/n \leq R_{\mathbb{E}}^*(\mathcal{U}(k)) \leq \sup_{\mu \in \mathcal{U}(k)} \mathbb{E}_\mu \|\hat{\mu}^{\text{LS}}(\mathcal{U}) - \mu\|_n^2 \leq c''\sigma^2 k \log(en/k)/n, \quad (6.3)$$

for some numerical constants $c', c'' > 0$ for instance by integrating the bounds of [3], Appendix C, where $\hat{\mu}^{\text{LS}}(\mathcal{U}) \in \operatorname{argmin}_{\mu \in \mathcal{U}} \|\mu - \mathbf{y}\|$ is the non-convex unimodal Least-Squares estimator (if the non-convex optimization problem has several solutions, we break ties arbitrarily).

The results below show that, using the number of constant pieces $\hat{k} = p(\hat{\mu}^{\text{LS}}(\mathcal{U}))$ of the unimodal Least-Squares estimator, one can construct adaptive confidence sets with respect to the collection of models (6.1).

Theorem 6.1. *Let $\hat{p} = k(\hat{\mu}^{\text{LS}}(\mathcal{U}))$ be the smallest integer such that $\hat{\mu}^{\text{LS}}(\mathcal{U})$ is piecewise constant on \hat{p} contiguous pieces. For all $\alpha \in (0, 1)$ and all $\mu \in \mathcal{U}$, the statistic*

$$\hat{r}_{\text{uni}} := \frac{4\sigma^2}{n} \left(\sqrt{(\hat{p}+3) \log\left(\frac{en}{\hat{p}+3}\right)} + \sqrt{2 \log\left(\frac{n^2}{\alpha}\right) + 2(\hat{p}+3) \log\left(\frac{en}{\hat{p}+3}\right)} \right)^2 \quad (6.4)$$

satisfies

$$\mathbb{P}_\mu \left(\|\hat{\mu}^{\text{LS}}(\mathcal{U}) - \mu\|_n^2 \leq \hat{r}_{\text{uni}} \right) \geq 1 - \alpha. \quad (6.5)$$

Proof. Let $\hat{\mu} = \hat{\mu}^{\text{LS}}(\mathcal{U})$ for brevity. By definition of $\hat{\mu}^{\text{LS}}(\mathcal{U})$, inequality $|\hat{\mu} - \mathbf{y}|_2^2 \leq |\mu - \mathbf{y}|_2^2$ holds. Hence, $|\hat{\mu} - \mu|_2^2 \leq 2\hat{\xi}^T(\hat{\mu} - \mu)$ and

$$|\hat{\mu} - \mu|_2 \leq 2\hat{\xi}^T \hat{\theta} \quad \text{where } \hat{\theta} = (\hat{\mu} - \mu)/|\hat{\mu} - \mu|_2. \quad (6.6)$$

We split $\{1, \dots, n\}$ into a partition (L, C, R) as follows: L is the largest set of indices of the form $\{1, 2, \dots, |L|\}$ such that both $\hat{\mu}_L$ and μ_L are non-increasing, R is the largest set of indices of the form $\{n - |R| + 1, n - |R| + 2, \dots, n\}$ such that $\hat{\mu}_R$ and μ_R are both non-decreasing, and $C = \{1, \dots, n\} \setminus (R \cup L)$ contains the remaining central indices, where by definition, $\hat{\mu}_C - \mu_C$ is either increasing or decreasing. Let \hat{k}_L be the number of constant pieces of $\hat{\mu}$ on L , and let $(\hat{S}_1, \dots, \hat{S}_{\hat{k}_L})$ be a partition of L such that $\hat{\mu}_L$ is constant on each \hat{S}_j . Similarly, let \hat{k}_R be the number of constant pieces of $\hat{\mu}$ on R and $(\hat{T}_1, \dots, \hat{T}_{\hat{k}_R})$ be a partition of R such that $\hat{\mu}_R$ is constant on each \hat{T}_j . Then both $\hat{\mu} - \mu$ and $\hat{\theta}$ belong to either one of the cones

$$\hat{\mathcal{D}}_+ := (\mathcal{S}_{|\hat{S}_1|}^\uparrow \times \dots \times \mathcal{S}_{|\hat{S}_{\hat{k}_L}|}^\uparrow) \times \mathcal{S}_{|C|}^\uparrow \times (\mathcal{S}_{|\hat{T}_1|}^\downarrow \times \dots \times \mathcal{S}_{|\hat{T}_{\hat{k}_R}|}^\downarrow).$$

or

$$\hat{\mathcal{D}}_- := (\mathcal{S}_{|\hat{S}_1|}^\uparrow \times \dots \times \mathcal{S}_{|\hat{S}_{\hat{k}_L}|}^\uparrow) \times \mathcal{S}_{|C|}^\downarrow \times (\mathcal{S}_{|\hat{T}_1|}^\downarrow \times \dots \times \mathcal{S}_{|\hat{T}_{\hat{k}_R}|}^\downarrow),$$

the only difference being the direction on the central indices in \mathcal{C} . We now argue similarly to (4.6) as follows. If \mathcal{C} is a nonrandom cone of the form

$$\mathcal{C} = \mathcal{S}_{n_1}^{\uparrow} \times \cdots \times \mathcal{S}_{n_k}^{\uparrow} \times \mathcal{S}_{m_1}^{\downarrow} \times \cdots \times \mathcal{S}_{m_l}^{\downarrow} \quad (6.7)$$

for positive integers $n_1, \dots, n_k, m_1, \dots, m_l$ such that $n_1 + \cdots + n_k + m_1 + \cdots + m_l = n$, then by the Gaussian concentration theorem, (4.5) holds with probability at least $1 - e^{-x}$. Furthermore, $\mathbb{E} \sup_{\theta \in \mathcal{C}: \|\theta\|_2=1} \xi^T \theta \leq \sigma \delta(\mathcal{C})^{1/2}$ by Jensen's inequality, and $\delta(\mathcal{C}) = \sum_{j=1}^k \log(en_j) + \sum_{j=1}^l \log(em_j) \leq (k+l) \log(en/(k+l))$ thanks to (3.12) and the fact that the statistical dimension of a product of cones is equal to the sum of the statistical dimensions (cf. [1], Proposition 3.1, or the discussion preceding (3.8)).

Let $s = 1, \dots, n$ be fixed. There are fewer than $n \binom{n-1}{s}$ cones of the form (6.7) with $k+l=s$. By the union bound and inequality $\log \binom{n-1}{s} \leq s \log(en/s)$, we have with probability at least $1 - e^{-x}$ the bound

$$\sup_{\mathcal{C}} \left(\sup_{\theta \in \mathcal{C}: \|\theta\|_2=1} \xi^T \theta \right) \leq \sigma \sqrt{s \log(en/s)} + \sigma \sqrt{2(x + \log n + s \log(en/s))}$$

where the supremum is taken over all cones \mathcal{C} of the form (6.7) with $k+l=s$. Finally, we apply the union bound over all $s \in \{1, \dots, n\}$ and set $x = \log(n/\alpha)$. We have established that with probability at least $1 - \alpha$,

$$\sup_{\theta \in \hat{\mathcal{D}}_+ \cup \hat{\mathcal{D}}_-: \|\theta\|_2=1} \xi^T \theta \leq \sigma \sqrt{\hat{s} \log(en/\hat{s})} + \sigma \sqrt{2(\log(n^2/\alpha) + \hat{s} \log(en/\hat{s}))}.$$

where $\hat{s} = \hat{k}_L + \hat{k}_R + 1$. The proof is completed by combining this bound with (6.6), and observing that $\hat{s} = \hat{k}_R + 1 + \hat{k}_L \leq \hat{p} + 3$. \square

Finally, the following result shows that the confidence set of the previous theorem has optimal radius up to logarithm factors.

Theorem 6.2. *If $\mu \in \mathcal{U}$, then*

$$\hat{p} := k(\hat{\mu}^{\text{LS}}(\mathcal{U})) \leq 4k(\mu) \log(en/k(\mu)) + 14 \log(2n/\gamma)$$

with probability at least $1 - \gamma$. Furthermore, $\mathbb{E}[\hat{p}] \leq 4k(\mu) \log(en/k(\mu)) + 14 \log(2en)$.

Proof. If $\hat{\mu} = \hat{\mu}^{\text{LS}}(\mathcal{U})$ is a unimodal fit and \hat{m} is the location of its mode, then the following facts hold true [28]: (a) the value of $\hat{\mu}_{\hat{m}}$ of $\hat{\mu}$ at \hat{m} is equal to $y_{\hat{m}}$, (b) $\hat{\mu}_{\{1, \dots, \hat{m}\}}$ is equal to the isotonic (decreasing) fit of $\mathbf{y}_{\{1, \dots, \hat{m}\}}$, and (c) $\hat{\mu}_{\{\hat{m}, \dots, n\}}$ is equal to the isotonic (increasing) fit of $\mathbf{y}_{\{\hat{m}, \dots, n\}}$. Hence, we can bound from above the number of constant pieces of $\hat{\mu}$ by the number of constant pieces of two isotonic fits, one on $\{1, \dots, \hat{m}\}$ and the other on $\{\hat{m}, \dots, n\}$.

If m is a deterministic mode location, and $\mu_{\{1, \dots, m\}}$ has $k(\mu_{\{1, \dots, m\}})$ pieces, by Theorem 4.2 the isotonic (decreasing) fit of $\mathbf{y}_{\{1, \dots, m\}}$ has at most

$$2k(\mu_{\{1, \dots, m\}}) \log(em/k(\mu_{\{1, \dots, m\}})) + 7 \log(1/\gamma)$$

constant pieces with probability $1 - \gamma$. Similarly, with at least probability $1 - \gamma$, the isotonic (increasing) fit of $\mathbf{y}_{\{m, \dots, n\}}$ has at most

$$2k(\mu_{\{m, \dots, n\}}) \log(e(n-m+1)/k(\mu_{\{m, \dots, n\}})) + 7 \log(1/\gamma)$$

constant pieces with probability at least $1 - \gamma$. By the union bound, the two previous sentences hold uniformly over all possible modes $m = 1, \dots, n$ with probability at least $1 - 2n\gamma$.

Hence, with probability at least $1 - 2n\gamma$, the number of constant pieces of the unimodal least-squares $\hat{\mu}$ is bounded from above by

$$2k(\mu_{\{1, \dots, \hat{m}\}}) \log\left(\frac{e\hat{m}}{k(\mu_{\{1, \dots, \hat{m}\}})}\right) + 2k(\mu_{\{\hat{m}, \dots, n\}}) \log\left(\frac{e(n - \hat{m} + 1)}{k(\mu_{\{\hat{m}, \dots, n\}})}\right) + 14 \log \frac{1}{\gamma}.$$

We first use that $\hat{m} \leq n$ and $(n - \hat{m} + 1) \leq n$ to bound from above the numerator inside the logarithms. Next, we use $k(\mu_{\{1, \dots, \hat{m}\}}) \leq k(\mu)$ and the fact that $x \log(en/x)$ is increasing on $[1, n]$ to conclude that the previous display is bounded from above by $4k(\mu) \log(en/k(\mu)) + 14 \log(1/\gamma)$. The result in expectation is obtained by integration, using the identity $\mathbb{E}[Z] = \int_0^\infty \mathbb{P}(Z > t) dt$ for every non-negative random variable Z . \square

Corollary 6.3. *Let $J = \{1, \dots, n\}$ and define the collection of models $(E_k)_{k \in J} = (\mathcal{U}(k))_{k=1, \dots, n}$. The random variable \hat{r}_{uni} defined in (6.4) satisfies (6.5), (2.8) and (2.9) with \hat{r} replaced by \hat{r}_{uni} . Thus, the ball centered at $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\cup})$ of radius $\sqrt{\hat{r}_{uni}}$ is an honest confidence set, which is adaptive in expectation with respect to the models $(\mathcal{U}(k))_{k=1, \dots, n}$.*

The proof follows by combining Theorems 6.1 and 6.2, similarly to the proof of Corollary 4.3 for non-decreasing sequences.

7. Concluding remarks

We have provided a simple construction of honest and adaptive confidence sets for isotonic, convex and unimodal regression. Our construction reveals that the complexity of the Least-Squares estimator in these problems, e.g. the number of jumps of the isotonic Least-Squares or the number of changes of slope of the convex Least-Squares, can be used to bound from above the error of the estimator (cf. Theorems 4.1, 5.2 and 6.1). Furthermore, the complexity of the Least-Squares estimator in these problems is not larger, up to logarithmic factors, than the complexity of the true mean vector (cf. Theorems 4.2, 5.4 and 6.2).

The construction of honest confidence sets in Theorems 4.1, 5.2 and 6.1 relies on a careful application of the Gaussian concentration theorem combined with union bounds and upper bounds on statistical dimensions of tangent cones. Such techniques can readily be extended to the setting of [12, 18, 21], where bounds on the statistical dimensions of tangent cones are readily available. However, the techniques used in Theorems 4.2, 5.4 and 6.2 to control the size of such confidence sets do not directly extend to the settings considered in [12, 18, 21] and it is unclear at this point how to control adaptively the radius of the confidence sets in these settings.

For the piecewise constant signals without the monotonicity constraint studied in Section 4.2, construction of honest and adaptive confidence sets are impossible. This provides additional evidence that shape constraints make it possible to construct adaptive and honest confidence sets or confidence bands for models where such task is impossible in the absence of the shape constraint [6, 7, 15].

The present paper tackles shape constraint regression problems with fixed design: if $x_1 < \dots < x_n$ are ordered design points then the model \mathcal{S}_n^{\uparrow} of nondecreasing sequences is equal to $\{(f(x_1), \dots, f(x_n))^T, f : \mathbb{R} \rightarrow \mathbb{R} \text{ nondecreasing}\}$. Hence, the results of the paper readily extend to random design isotonic regression, provided that the loss is measured with the empirical loss $\frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f^*(x_i))^2$ where f^* is the true regression function and $x_1 < \dots < x_n$ are the random design points, by setting $\mu_i = f^*(x_i)$. For convex regression, if $x_1 < \dots < x_n$ are equispaced

(i.e., $|x_{i+1} - x_i|$ is the same for all i), the model \mathcal{S}_n^{\cup} is the same as $\{(f(x_1), \dots, f(x_n))^T, f: \mathbb{R} \rightarrow \mathbb{R} \text{ convex}\}$. If the design points are not equispaced, then $\{(f(x_1), \dots, f(x_n))^T, f: \mathbb{R} \rightarrow \mathbb{R} \text{ convex}\}$ defines a model different than \mathcal{S}_n^{\cup} , but the same logarithmic bounds on the statistical dimension hold for non-equispaced design points [3], cf. Proposition 4.2, and consequently the results of Section 5 extend to non-equispaced designs.

Regarding random design, in order to extend the results of the present paper to i.i.d. real valued design points x_1, \dots, x_n , one would need to establish a relationship between the empirical loss $\frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f^*(x_i))^2$ and the population loss $\mathbb{E}_{x \sim F}[(\hat{f}(x) - f^*(x))^2 | \hat{f}]$ where F is the distribution of the design. The random design analysis in [21] is a good starting point for this purpose. Finally, let us mention that if the design is Gaussian and high-dimensional, upper bounds based on statistical dimensions can be obtained similarly as in the sequence model [23] and some techniques of the present paper may thus be applicable. We leave this research direction open for future work.

Appendix A: Nondecreasing sequences with bounded total variation

Let $V > 0$. If the unknown parameter μ satisfies $\mu_n - \mu_1 \leq V$, the risk of the Least Squares estimator satisfy [31], (28),

$$\mathbb{E}_{\mu} \|\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\uparrow}) - \mu\|_n^2 \leq \sigma^2 \kappa^2 \left(\left(\frac{V}{\sigma n} \right)^{2/3} + \frac{\log(en)}{n} \right),$$

where $\kappa \leq 3.6$. Thus, an explicit constant is readily available [31], (2.8).

In this section, we explain how to construct confidence sets with diameter of the same order as the right-hand side of the previous display. We proceed as follows.

The function $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ defined by $f(\mathbf{v}) = \|\Pi_{\mathcal{S}_n^{\uparrow}}(\mu + \sigma \mathbf{v}) - \mu\|_n$ is Lipschitz with coefficient σ/\sqrt{n} as for all $\mathbf{v}, \mathbf{v}' \in \mathbb{R}^n$,

$$\begin{aligned} |f(\mathbf{v}) - f(\mathbf{v}')| &\leq \|\Pi_{\mathcal{S}_n^{\uparrow}}(\mu + \sigma \mathbf{v}) - \Pi_{\mathcal{S}_n^{\uparrow}}(\mu + \sigma \mathbf{v}')\|_n \\ &\leq \sigma \|\mathbf{v} - \mathbf{v}'\|_n = (\sigma/\sqrt{n}) |\mathbf{v} - \mathbf{v}'|_2. \end{aligned} \tag{A.1}$$

By the Gaussian concentration inequality [5], Theorem 5.6, the following holds with probability greater than $1 - \alpha$

$$\|\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\uparrow}) - \mu\|_n \leq \mathbb{E}_{\mu} \|\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\uparrow}) - \mu\|_n + \sigma \sqrt{\frac{2 \log(1/\alpha)}{n}}.$$

Using that $(a+b)^2 \leq 2a^2 + 2b^2$, we obtain the following for all $\alpha \in (0, 1)$: If $\mu \in \mathcal{S}_n^{\uparrow}$ and $\mu_n - \mu_1 \leq V$, then

$$\|\hat{\mu}^{\text{LS}}(\mathcal{S}_n^{\uparrow}) - \mu\|_n^2 \leq 2\kappa^2 \sigma^2 \left(\frac{V}{\sigma n} \right)^{2/3} + \frac{2\kappa^2 \sigma^2 \log(en) + 4\sigma^2 \log(1/\alpha)}{n}$$

with probability greater than $1 - \alpha$.

Let $V_{\mu} = \mu_n - \mu_1$ and $\hat{V} = y_n - y_1$. The random variable $\hat{V} - V_{\mu}$ is a centered Gaussian with variance $2\sigma^2$, so that

$$V_{\mu} \leq \hat{V} + 2\sigma \sqrt{\log(1/\alpha)}$$

with probability greater than $1 - \alpha$. Thus, we have established the following.

Proposition A.1. Let $\mu \in \mathcal{S}_n^\uparrow$. Define the statistic \hat{s}_\uparrow by

$$\hat{s}_\uparrow = 2\kappa^2\sigma^2 \left(\frac{\hat{V} + 2\sigma\sqrt{\log(1/\alpha)}}{\sigma n} \right)^{2/3} + \frac{2\kappa^2\sigma^2 \log(en) + 4\sigma^2 \log(1/\alpha)}{n}$$

where $\kappa \leq 3.6$ is the constant from [31] that appears in (A.1). Then we have $\|\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow) - \mu\|_n^2 \leq \hat{s}_\uparrow$ with probability greater than $1 - 2\alpha$.

Furthermore, it is clear that $\hat{V} \leq V_\mu + 2\sigma\sqrt{\log(1/\gamma)}$ with probability greater than $1 - \gamma$ for all $\gamma \in (0, 1)$.

Proposition A.2. Let $\mu \in \mathcal{S}_n^\uparrow$ and let $V = \mu_n - \mu_1$. Then the statistic \hat{s}_\uparrow defined above satisfies

$$\hat{s}_\uparrow \leq 2\kappa^2\sigma^2 \left(\frac{\hat{V} + 2\sigma\sqrt{\log(1/(\gamma\alpha))}}{\sigma n} \right)^{2/3} + \frac{2\kappa^2\sigma^2 \log(en) + 4\sigma^2 \log(1/\alpha)}{n} \quad (\text{A.2})$$

with probability at least $1 - \gamma$ for all $\gamma \in (0, 1)$.

Theorem A.3. Let $\mu \in \mathcal{S}_n^\uparrow$. The statistic $\min(\hat{r}_\uparrow, \hat{s}_\uparrow)$ satisfies

$$\|\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\uparrow) - \mu\|_n^2 \leq \min(\hat{r}_\uparrow, \hat{s}_\uparrow)$$

with probability at least $1 - 3\alpha$. Furthermore, for all $\gamma \in (0, 1)$, the statistic $\min(\hat{r}_\uparrow, \hat{s}_\uparrow)$ is bounded from above with probability at least $1 - 2\gamma$, by the minimum of the right-hand side of (4.7) and the right-hand side of (A.2).

For all $V \geq \sigma$ and all $k = 1, \dots, n$, define the class

$$\mathcal{S}_n^\uparrow(k, V) := \{v = (v_1, \dots, v_n)^T \in \mathcal{S}_n^\uparrow : k(v) \leq k \text{ and } v_n - v_1 \leq V\}.$$

Since

$$c\sigma^2 \min\left(\left(\frac{V}{\sigma n}\right)^{2/3}, \frac{k \log \log(16n/k)}{n}\right) \leq R_{\mathbb{E}}^*(\mathcal{S}_n^\uparrow(k, V))$$

for some absolute constant $c > 0$ by combining the lower bound in (3.2) with the lower bound in [31], (2.9). Thus, the statistic $\min(\hat{r}_\uparrow, \hat{s}_\uparrow)$ of Theorem A.3 provides an honest confidence ball centered at the Least Squares estimator, and this confidence ball is adaptive in expectation for the collection of models

$$(\mathcal{S}_n^\uparrow(k, V))_{k \in \{1, \dots, n\}, V \in [\sigma, +\infty)}.$$

Appendix B: Technical lemma

Lemma B.1. In the present lemma, all quantities are deterministic. Let $a \in [-\infty, +\infty)$ and $b \in (-\infty, +\infty]$ such that $a \leq b$. Let $\mathbf{y} \in \mathbb{R}^n$. Define $\boldsymbol{\theta}$ and $\boldsymbol{\theta}^*$ as the unique solutions of the minimization problems

$$\boldsymbol{\theta}^* \in \operatorname{argmin}_{\mathbf{v} \in \mathcal{S}_n^\uparrow} \|\mathbf{y} - \mathbf{v}\|_2^2, \quad (\text{B.1})$$

$$\boldsymbol{\theta} \in \underset{\mathbf{v} \in \mathcal{S}_n^{\uparrow}(a, b)}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{v}\|_2^2 \quad (\text{B.2})$$

where $\mathcal{S}_n^{\uparrow}(a, b) := \{\mathbf{v} = (v_1, \dots, v_n)^T \in \mathcal{S}_n^{\uparrow} : a \leq v_1, v_n \leq b\}$. Then $k(\boldsymbol{\theta}) \leq k(\boldsymbol{\theta}^*)$.

The intuition behind this lemma is that if a constraint is not saturated for $\boldsymbol{\theta}$, this constraint is not saturated for $\boldsymbol{\theta}^*$ either, so $\boldsymbol{\theta}^*$ has at least as many jumps as $\boldsymbol{\theta}$.

Proof of Lemma B.1. Let $T_a = \{i = 1, \dots, n : \hat{\theta}_i^* \leq a\}$, $T_b = \{i = 1, \dots, n : \hat{\theta}_i^* \geq b\}$ and $T_c = \{i = 1, \dots, n : a < \hat{\theta}_i^* < b\}$. We will prove that the unique minimizer $\boldsymbol{\theta}$ of the problem (B.2) is

$$\boldsymbol{\theta}_{T_a} = a\mathbf{1}_{T_a}, \quad \boldsymbol{\theta}_{T_c} = \boldsymbol{\theta}_{T_c}^*, \quad \boldsymbol{\theta}_{T_b} = b\mathbf{1}_{T_b}, \quad (\text{B.3})$$

where $\mathbf{1} = (1, \dots, 1)^T \in \mathbb{R}^n$. Then it is clear that $k(\boldsymbol{\theta}) = 1 + k(\boldsymbol{\theta}_{T_c}^*) + 1 \leq k(\boldsymbol{\theta}_{T_a}^*) + k(\boldsymbol{\theta}_{T_c}^*) + k(\boldsymbol{\theta}_{T_b}^*) = k(\boldsymbol{\theta}^*)$.

First, by strong convexity there exists a unique solution to the minimization problem (B.2), and a unique solution to the minimization problem (B.1). Second, by the characterization of the projection onto the closed convex set $\mathcal{S}_n^{\uparrow}(a, b)$, if $\boldsymbol{\theta}$ satisfies

$$A_{\mathbf{u}} := (\mathbf{u} - \boldsymbol{\theta})^T (\mathbf{y} - \boldsymbol{\theta}) \leq 0$$

for all $\mathbf{u} \in \mathcal{S}_n^{\uparrow}(a, b)$, then $\boldsymbol{\theta}$ is the unique solution to the minimization problem (B.2). Let $\boldsymbol{\theta}$ be defined by (B.3). By simple algebra, for all $\mathbf{u} \in \mathcal{S}_n^{\uparrow}(a, b)$,

$$\begin{aligned} A_{\mathbf{u}} &= (\mathbf{u}_{T_a} - a\mathbf{1}_{T_a} + \boldsymbol{\theta}_{T_a}^* - \boldsymbol{\theta}_{T_a}^*)^T (\mathbf{y}_{T_a} - \boldsymbol{\theta}_{T_a}^*) + (\mathbf{u}_{T_a} - a\mathbf{1}_{T_a})^T (\boldsymbol{\theta}_{T_a}^* - a\mathbf{1}_{T_a}) \\ &\quad + (\mathbf{u}_{T_b} - b\mathbf{1}_{T_b} + \boldsymbol{\theta}_{T_b}^* - \boldsymbol{\theta}_{T_b}^*)^T (\mathbf{y}_{T_b} - \boldsymbol{\theta}_{T_b}^*) + (\mathbf{u}_{T_b} - b\mathbf{1}_{T_b})^T (\boldsymbol{\theta}_{T_b}^* - b\mathbf{1}_{T_b}) \\ &\quad + (\mathbf{u}_{T_c} - \boldsymbol{\theta}_{T_c}^*)^T (\mathbf{y}_{T_c} - \boldsymbol{\theta}_{T_c}^*). \end{aligned}$$

If a vector \mathbf{v} has nonnegative entries and a vector \mathbf{x} has non-positive entries, then $\mathbf{v}^T \mathbf{x} \leq 0$, so $(\mathbf{u}_{T_a} - a\mathbf{1}_{T_a})^T (\boldsymbol{\theta}_{T_a}^* - a\mathbf{1}_{T_a}) \leq 0$ and $(\mathbf{u}_{T_b} - b\mathbf{1}_{T_b})^T (\boldsymbol{\theta}_{T_b}^* - b\mathbf{1}_{T_b}) \leq 0$. Thus,

$$\begin{aligned} A_{\mathbf{u}} &\leq (\mathbf{u}_{T_a} - a\mathbf{1}_{T_a} + \boldsymbol{\theta}_{T_a}^* - \boldsymbol{\theta}_{T_a}^*)^T (\mathbf{y}_{T_a} - \boldsymbol{\theta}_{T_a}^*) \\ &\quad + (\mathbf{u}_{T_b} - b\mathbf{1}_{T_b} + \boldsymbol{\theta}_{T_b}^* - \boldsymbol{\theta}_{T_b}^*)^T (\mathbf{y}_{T_b} - \boldsymbol{\theta}_{T_b}^*) \\ &\quad + (\mathbf{u}_{T_c} - \boldsymbol{\theta}_{T_c}^*)^T (\mathbf{y}_{T_c} - \boldsymbol{\theta}_{T_c}^*), \end{aligned}$$

and the right-hand side of the previous display is equal to

$$(\mathbf{v} - \boldsymbol{\theta}^*)^T (\mathbf{y} - \boldsymbol{\theta}^*), \quad (\text{B.4})$$

where $\mathbf{v} \in \mathbb{R}^n$ is defined by

$$\mathbf{v}_{T_a} := \mathbf{u}_{T_a} - a\mathbf{1}_{T_a} + \boldsymbol{\theta}_{T_a}^*, \quad \mathbf{v}_{T_c} := \mathbf{u}_{T_c}, \quad \mathbf{v}_{T_b} := \mathbf{u}_{T_b} - b\mathbf{1}_{T_b} + \boldsymbol{\theta}_{T_b}^*.$$

We have $\mathbf{v} \in \mathcal{S}_n^{\uparrow}$ by definition of T_a , T_c and T_b . The quantity (B.4) is non-positive because $\boldsymbol{\theta}^*$ is the projection of \mathbf{y} onto the convex set \mathcal{S}_n^{\uparrow} . Thus, we have established that $A_{\mathbf{u}} \leq 0$ for all $\mathbf{u} \in \mathcal{S}_n^{\uparrow}(a, b)$, so that the unique solution of the minimization problem (B.2) is given by the expression (B.3). \square

Appendix C: Proofs for convex sequences

Proof of Theorem 5.2. For any set T of the form (4.3), using the concentration property (3.7) of the random variable (3.6) with $K = \mathcal{S}_{|T|}^\cap$, we have with probability greater than $1 - \alpha$,

$$|\Pi_{\mathcal{S}_{|T|}^\cap}(\xi_T)|_2^2 \leq 2\delta(\mathcal{S}_{|T|}^\cap) + 4\log(1/\alpha) \leq 16\log(en) + 4\log(1/\alpha),$$

where we used (5.2) for the last inequality. There are less than n^2 sets $T \subset \{1, \dots, n\}$ of the form (4.3). Using the union bound for all sets T of the form (4.3), we have $\mathbb{P}(\Omega(\alpha)) \geq 1 - \alpha$ where

$$\Omega(\alpha) := \left\{ \sup_{T \subset \{1, \dots, n\} \text{ of the form (4.3)}} |\Pi_{\mathcal{S}_{|T|}^\cap}(\xi_T)|_2^2 \leq \sigma^2 \left(16\log(en) + 4\log\left(\frac{n^2}{\alpha}\right) \right) \right\}.$$

Let $\hat{\mu} = \hat{\mu}^{\text{LS}}(\mathcal{S}_n^\cup)$ for notational simplicity. Then (2.6) with μ replaced by μ can be rewritten as

$$|\hat{\mu} - \mu|_2^2 \leq 2\xi^T(\hat{\mu} - \mu) - |\hat{\mu} - \mu|_2^2.$$

By definition of $q(\cdot)$, there exists a partition $(\hat{T}_1, \dots, \hat{T}_{\hat{q}})$ of $\{1, \dots, n\}$ such that $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\cup)$ is affine on each \hat{T}_j , $j = 1, \dots, \hat{q}$. Furthermore, each \hat{T}_j has the form (4.3) because $\hat{\mu}^{\text{LS}}(\mathcal{S}_n^\cup) \in \mathcal{S}_n^\cup$. We have

$$\begin{aligned} 2\xi^T(\hat{\mu} - \mu) - |\hat{\mu} - \mu|_2^2 &= \sum_{j=1}^{\hat{q}} 2\xi_{\hat{T}_j}^T(\hat{\mu} - \mu)_{\hat{T}_j} - |(\hat{\mu} - \mu)_{\hat{T}_j}|_2^2, \\ &\leq \sum_{j=1}^{\hat{q}} \left(\frac{\xi_{\hat{T}_j}^T(\hat{\mu} - \mu)_{\hat{T}_j}}{|(\hat{\mu} - \mu)_{\hat{T}_j}|_2} \right)^2, \end{aligned}$$

where we have used $2ab - a^2 \leq b^2$. By definition of $(\hat{T}_1, \dots, \hat{T}_{\hat{q}})$, $\hat{\mu}$ is affine on \hat{T}_j for each $j = 1, \dots, \hat{q}$, thus the vector $(\hat{\mu} - \mu)_{\hat{T}_j} \in \mathcal{S}_{|\hat{T}_j|}^\cap$ is a concave sequence. By taking the supremum, we obtain

$$|\hat{\mu} - \mu|_2^2 \leq \sum_{j=1}^{\hat{q}} \sup_{\mathbf{v} \in \mathcal{S}_{|\hat{T}_j|}^\cap : |\mathbf{v}|_2^2 \leq 1} (\xi_{\hat{T}_j}^T \mathbf{v})^2 = \sum_{j=1}^{\hat{q}} |\Pi_{\mathcal{S}_{|\hat{T}_j|}^\cap}(\xi_{|\hat{T}_j|})|_2^2,$$

where we used (3.6) for the last equality. On the event $\Omega(\alpha)$ and by definition of \hat{r}_\cup ,

$$|\hat{\mu} - \mu|_2^2 \leq \sum_{j=1}^{\hat{q}} |\Pi_{\mathcal{S}_{|\hat{T}_j|}^\cap}(\xi_{|\hat{T}_j|})|_2^2 \leq \sigma^2 \hat{q} (16\log(en) + 4\log(n^2/\alpha)) = n\hat{r}_\cup. \quad \square$$

Acknowledgements

This work was partially supported by GENES and by the ANR grant ANR-11-IDEX-0003/Labex Ecodec/ANR-11-LABX-0047, as well as the NSF grants DMS-1811976 and DMS-1945428.

References

- [1] Amelunxen, D., Lotz, M., McCoy, M.B. and Tropp, J.A. (2014). Living on the edge: Phase transitions in convex programs with random data. *Inf. Inference* **3** 224–294. [MR3311453](#) <https://doi.org/10.1093/imaiai/iau005>
- [2] Baraud, Y. (2002). Non-asymptotic minimax rates of testing in signal detection. *Bernoulli* **8** 577–606. [MR1935648](#)
- [3] Bellec, P.C. (2018). Sharp oracle inequalities for least squares estimators in shape restricted regression. *Ann. Statist.* **46** 745–780. [MR3782383](#) <https://doi.org/10.1214/17-AOS1566>
- [4] Bellec, P.C. and Tsybakov, A.B. (2015). Sharp oracle bounds for monotone and convex regression through aggregation. *J. Mach. Learn. Res.* **16** 1879–1892. [MR3417801](#)
- [5] Boucheron, S., Lugosi, G. and Massart, P. (2013). *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford: Oxford Univ. Press. [MR3185193](#) <https://doi.org/10.1093/acprof:oso/9780199535255.001.0001>
- [6] Cai, T.T. and Low, M.G. (2006). Adaptive confidence balls. *Ann. Statist.* **34** 202–228. [MR2275240](#) <https://doi.org/10.1214/009053606000000146>
- [7] Cai, T.T., Low, M.G. and Xia, Y. (2013). Adaptive confidence intervals for regression functions under shape constraints. *Ann. Statist.* **41** 722–750. [MR3099119](#) <https://doi.org/10.1214/12-AOS1068>
- [8] Carpentier, A. and Verzelen, N. (2019). Optimal sparsity testing in linear regression model. ArXiv preprint. Available at [arXiv:1901.08802](#).
- [9] Carpentier, A. and Verzelen, N. (2019). Adaptive estimation of the sparsity in the Gaussian vector model. *Ann. Statist.* **47** 93–126. [MR3909928](#) <https://doi.org/10.1214/17-AOS1680>
- [10] Chandrasekaran, V., Recht, B., Parrilo, P.A. and Willsky, A.S. (2012). The convex geometry of linear inverse problems. *Found. Comput. Math.* **12** 805–849. [MR2989474](#) <https://doi.org/10.1007/s10208-012-9135-7>
- [11] Chatterjee, S., Guntuboyina, A. and Sen, B. (2015). On risk bounds in isotonic and other shape restricted regression problems. *Ann. Statist.* **43** 1774–1800. [MR3357878](#) <https://doi.org/10.1214/15-AOS1324>
- [12] Chatterjee, S., Guntuboyina, A. and Sen, B. (2018). On matrix estimation under monotonicity constraints. *Bernoulli* **24** 1072–1100. [MR3706788](#) <https://doi.org/10.3150/16-BEJ865>
- [13] Chatterjee, S. and Lafferty, J. (2019). Adaptive risk bounds in unimodal regression. *Bernoulli* **25** 1–25. [MR3892309](#) <https://doi.org/10.3150/16-bej922>
- [14] Chen, X., Lin, Q. and Sen, B. (2020). On Degrees of Freedom of Projection Estimators With Applications to Multivariate Nonparametric Regression. *J. Amer. Statist. Assoc.* **115** 173–186. [MR4078455](#) <https://doi.org/10.1080/01621459.2018.1537917>
- [15] Dümbgen, L. (2003). Optimal confidence bands for shape-restricted curves. *Bernoulli* **9** 423–449. [MR1997491](#) <https://doi.org/10.3150/bj/1065444812>
- [16] Flammarion, N., Mao, C. and Rigollet, P. (2019). Optimal rates of statistical seriation. *Bernoulli* **25** 623–653. [MR3892331](#) <https://doi.org/10.3150/17-bej1000>
- [17] Gao, C., Han, F. and Zhang, C.-H. (2017). Minimax risk bounds for piecewise constant models. ArXiv preprint. Available at [arXiv:1705.06386](#).
- [18] Guntuboyina, A., Lieu, D., Chatterjee, S. and Sen, B. (2020). Adaptive risk bounds in univariate total variation denoising and trend filtering. *Ann. Statist.* **48** 205–229. [MR4065159](#) <https://doi.org/10.1214/18-AOS1799>
- [19] Guntuboyina, A. and Sen, B. (2015). Global risk bounds and adaptation in univariate convex regression. *Probab. Theory Related Fields* **163** 379–411. [MR3405621](#) <https://doi.org/10.1007/s00440-014-0595-3>
- [20] Guntuboyina, A. and Sen, B. (2018). Nonparametric shape-restricted regression. *Statist. Sci.* **33** 568–594. [MR3881209](#) <https://doi.org/10.1214/18-STS665>
- [21] Han, Q., Wang, T., Chatterjee, S. and Samworth, R.J. (2019). Isotonic regression in general dimensions. *Ann. Statist.* **47** 2440–2471. [MR3988762](#) <https://doi.org/10.1214/18-AOS1753>
- [22] Ingster, Y.I., Tsybakov, A.B. and Verzelen, N. (2010). Detection boundary in sparse regression. *Electron. J. Stat.* **4** 1476–1526. [MR2747131](#) <https://doi.org/10.1214/10-EJS589>
- [23] Matey, N. (2019). Gaussian regression with convex constraints. In *Proceedings of Machine Learning Research* (K. Chaudhuri and M. Sugiyama, eds.). *Proceedings of Machine Learning Research* **89** 31–38. Available at <http://proceedings.mlr.press/v89/neykov19b.html>.

- [24] McCoy, M.B. and Tropp, J.A. (2014). From Steiner formulas for cones to concentration of intrinsic volumes. *Discrete Comput. Geom.* **51** 926–963. [MR3216671](#) <https://doi.org/10.1007/s00454-014-9595-4>
- [25] Meyer, M. and Woodroofe, M. (2000). On the degrees of freedom in shape-restricted regression. *Ann. Statist.* **28** 1083–1104. [MR1810920](#) <https://doi.org/10.1214/aos/1015956708>
- [26] Nickl, R. and van de Geer, S. (2013). Confidence sets in sparse regression. *Ann. Statist.* **41** 2852–2876. [MR3161450](#) <https://doi.org/10.1214/13-AOS1170>
- [27] Oymak, S. and Hassibi, B. (2016). Sharp MSE bounds for proximal denoising. *Found. Comput. Math.* **16** 965–1029. [MR3529131](#) <https://doi.org/10.1007/s10208-015-9278-4>
- [28] Stout, Q.F. (2008). Unimodal regression via prefix isotonic regression. *Comput. Statist. Data Anal.* **53** 289–297. [MR2649085](#) <https://doi.org/10.1016/j.csda.2008.08.005>
- [29] Tsybakov, A.B. (2009). *Introduction to Nonparametric Estimation. Springer Series in Statistics*. New York: Springer. [MR2724359](#) <https://doi.org/10.1007/b13794>
- [30] Verzelen, N. (2012). Minimax risks for sparse regressions: Ultra-high dimensional phenomenons. *Electron. J. Stat.* **6** 38–90. [MR2879672](#) <https://doi.org/10.1214/12-EJS666>
- [31] Zhang, C.-H. (2002). Risk bounds in isotonic regression. *Ann. Statist.* **30** 528–555. [MR1902898](#) <https://doi.org/10.1214/aos/1021379864>

Received April 2019 and revised March 2020