

Approximations of Semicontinuous Functions with Applications to Stochastic Optimization and Statistical Estimation

Johannes O. Royset
Operations Research Department
Naval Postgraduate School
joroyset@nps.edu

Abstract. Upper semicontinuous (usc) functions arise in the analysis of maximization problems, distributionally robust optimization, and function identification, which includes many problems of non-parametric statistics. We establish that every usc function is the limit of a hypo-converging sequence of piecewise affine functions of the difference-of-max type and illustrate resulting algorithmic possibilities in the context of approximate solution of infinite-dimensional optimization problems. In an effort to quantify the ease with which classes of usc functions can be approximated by finite collections, we provide upper and lower bounds on covering numbers for bounded sets of usc functions under the Attouch-Wets distance. The result is applied in the context of stochastic optimization problems defined over spaces of usc functions. We establish confidence regions for optimal solutions based on sample average approximations and examine the accompanying rates of convergence. Examples from nonparametric statistics illustrate the results.

Keywords: hypo-convergence, Attouch-Wets distance, approximation theory, solution stability, stochastic optimization, epi-splines, rate of convergence.

Date: July 9, 2019

1 Introduction

Extended real-valued upper-semicontinuous (usc) functions on \mathbb{R}^n are fundamental in the study of finite-dimensional constrained maximization problems as essentially all such problems can be represented by usc functions. They arise in probability theory with distribution and càdlàg functions also being usc. Emerging applications of usc functions in infinite-dimensional problems include nonparametric statistical M -estimation [38], distributionally robust optimization [36], and more generally function identification [35]. In these applications, optimization problems are formulated over spaces of usc functions. Regardless of the setting, it becomes important to have means to approximate usc functions as well as an understanding of the difficulty with such an undertaking. This article provides three main results in these directions: (i) We establish that every usc function is the limit of a hypo-converging sequence of mesh-free piecewise affine functions of the difference-of-max type. Thus, as a corollary, the difference-of-convex (dc) functions are hypo-dense in spaces of usc functions. With the advances

in computational treatment of dc functions (see for example [8]), this leads to numerous algorithmic possibilities, which we illustrate in the context of function identification problems. (ii) We provide upper and lower bounds on covering numbers for bounded sets of usc functions under the Attouch-Wets (aw) distance and thereby quantify the ease with which classes of usc functions can be approximated by finite collections. (iii) For stochastic optimization problems defined over spaces of usc functions, we establish confidence regions for optimal solutions in terms of the aw-distance and sample average approximations, with rates of convergence as the sample size grows. The result requires only semicontinuity of the objective function and therefore applies in challenging settings such as simulation optimization of “black-box” stochastic systems where little structure may be known.

The consideration of approximations in the sense of hypo-convergence, which is metrized by the aw-distance, is natural and convenient in many applications. If an usc function is approximated in this sense, then the maximizers of the approximating function will be “near” those of the actual function. This is exactly the desired property when the usc function represents a constrained maximization problem. It is also the goal when the usc function is a probability density function and we need to estimate its modes; the approximating density will have modes “near” the actual modes. The situation is similar when the usc function is a surrogate model in an engineering design problem that needs to be maximized to find an optimal design; see Section 5 for an example. The notion of approximation is further motivated in the context of distribution functions by the fact that for such functions hypo-convergence is equivalent to convergence in distribution, a property that is leveraged to address optimization under stochastic ambiguity in [36]. An alternative focus on approximations in the sense of uniform convergence would have limited the scope to finite-valued continuous functions with common compact domains, which is too restrictive in many applications. Hypo-convergence permits treatment of usc functions defined on any subset of \mathbb{R}^n .

The study of hypo-converging usc functions and, in parallel, epi-converging extended real-valued lower-semicontinuous (lsc) functions has a long history, with important accomplishments in convex and nonsmooth analysis as well as the approximation theory of maximization and minimization problems; see [31] for details. Connections to probability theory are established in [39, 40] and more recently in [36]. The first formulation of infinite-dimensional optimization problems over spaces of semicontinuous functions appears in [34], with theoretical developments in [35]. In particular, the latter reference defines the class of epi-splines (see also [33]), which are piecewise polynomial functions, and establishes that the class is dense in spaces of semicontinuous functions under the aw-distance. Even though epi-splines furnish a means to approximate arbitrary semicontinuous functions using a finite number of parameters, they suffer from the need to partition \mathbb{R}^n into a finite number of subsets. In the present paper, we show that semicontinuous functions can be approximated by piecewise affine functions that are defined *without specifying a partition* and that are characterized structurally as being the difference of two functions of the form $x \mapsto \max_{k=1,\dots,p} \langle a^k, x \rangle + \alpha_k$. Consequently, we refer to these piecewise affine functions as *mesh free*; the domain of each affine component adapts and is not preselected. This is a significant feature for medium- and high-dimensional problems, where representative low-dimensional subspaces need to be discovered and exploited and standard polynomial approximations become challenging (see [47] for some progress in such directions). Our approximation result for usc functions extends the well-known

fact that every continuous function on a convex compact set is the limit of a uniformly convergent sequence of dc functions, which can be traced back to the local property of dc functions established by [18]; see for example Proposition 2.3 in [21].

Covering numbers express the size of a class of functions in terms of the smallest number of balls with a certain radius needed to cover the class and are central to most consistency, rate of convergence, and error analysis in (non)parametric estimation and machine learning; see for example [44, 45, 16]. The pioneering work [23, 5] deal with continuous and smooth functions; see [30] for a more recent discussion. Functions of bounded variation are considered in [3] and analytic functions in [7]. An upper estimate for the covering numbers of the unit ball of Gaussian reproducing kernel Hilbert spaces is given in [48], with further refinements and applications in [46, 24]. Covering numbers of sets of convex functions are established in [11, 6], with significant improvements in [14]. The present paper establishes an upper bound on the covering numbers of bounded classes of usc functions under the aw-distance and show that it is sharp within a logarithmic factor.

Although sample average approximations are often used to solving stochastic optimization problems, it remains challenging to assess the quality of a solution obtained through such approximations. Upper and lower bounds on minimum values can be computed using the approaches in [20, 29, 27, 4] (see also [42, Sect. 5.6]), at least when problem relaxations can be solved to near global optimality. Validation approaches based on optimality conditions are found in [19, 43, 32, 26, 25] and [42, Sect. 5.6]. Rates of convergence of optimization problems with Lipschitz continuous objective functions defined on a compact subset of \mathbb{R}^n are given in [42, Sect. 5.3]. We leverage the results on covering numbers to establish confidence regions of optimal solutions of infinite-dimensional stochastic optimization problems defined on spaces of usc functions without assuming Lipschitz continuity. The result is novel even when specialized to finite dimensions. For a Hölder continuous case, we obtain, in some sense, a stronger result.

After a section laying out notation and terminology, we proceed in Section 3 with the result on piecewise affine approximations and its applications. Section 4 establishes bounds on covering numbers. Section 5 constructs confidence regions and discusses rates of convergence for solutions of stochastic optimization problems and their applications to nonparametric estimation. The paper ends with an appendix supplementing a proof.

2 Preliminaries

In some applications, it would be natural and beneficial to consider usc functions defined only on a strict subset of \mathbb{R}^n and their extensions to the whole \mathbb{R}^n by assigning the value $-\infty$ may not be meaningful. For example, if an usc function represents a necessarily nonnegative probability density, then such an assignment would not imply a useful extension. Consequently, we develop most results for usc functions defined on a nonempty closed subset $S \subset \mathbb{R}^n$, which could be all of \mathbb{R}^n , and is assumed to include the origin. Throughout, S will be such a set and the analysis will usually take place on the metric spaces $(S, \|\cdot - \cdot\|_\infty)$ and $(S \times \mathbb{R}, \|\cdot - \cdot\|_\infty)$; the difference from the usual $(\mathbb{R}^n, \|\cdot - \cdot\|_\infty)$ is anyhow minor and will be highlighted when significant. Of course, the sup-norm can be replaced by any other norm, but

this choice simplifies some expressions in Section 4. Likewise, the assumption $0 \in S$ can be relaxed, with the introduction of additional notation better avoided here. The facts of this section can be found in or deduced from [31, Chapter 7] and [33].

We recall that $\text{hypo } f = \{(x, \alpha) \in S \times \mathbb{R} \mid f(x) \geq \alpha\} \subset S \times \mathbb{R}$ is the *hypograph* of a function $f : S \rightarrow \overline{\mathbb{R}} = [-\infty, \infty]$. The collection of *usc functions* on S is denoted by

$$\text{usc-fcns}(S) = \{f : S \rightarrow \overline{\mathbb{R}} \mid \text{hypo } f \text{ is nonempty and closed}\}.$$

We let $\mathbb{N} = \{1, 2, \dots\}$. The *outer limit* of a sequence of sets $\{A^\nu, \nu \in \mathbb{N}\}$ in a topological space, denoted by $\text{OutLim } A^\nu$, is the collection of points to which a subsequence of $\{a^\nu \in A^\nu, \nu \in \mathbb{N}\}$ converges. The *inner limit*, denoted by $\text{InnLim } A^\nu$, is the collection of points to which a sequence $\{a^\nu \in A^\nu, \nu \in \mathbb{N}\}$ converges. If both limits exist and are equal to A , we say that $\{A^\nu, \nu \in \mathbb{N}\}$ *set-converges* to A and write $A^\nu \rightarrow A$ or $\text{Lim } A^\nu = A$. We denote by $\text{int } A$ and $\text{cl } A$ the interior and closure of A , respectively.

For $f^\nu, f \in \text{usc-fcns}(S)$,

$$f^\nu \text{ hypo-converges to } f, \text{ written } f^\nu \rightarrow f \iff \text{hypo } f^\nu \rightarrow \text{hypo } f.$$

Set-convergence of hypographs in this case, and therefore also hypo-convergence, is equivalent to having

$$\forall x^\nu \in S \rightarrow x, \quad \limsup f^\nu(x^\nu) \leq f(x) \tag{1}$$

$$\forall x \in S, \exists x^\nu \in S \rightarrow x \text{ with } \liminf f^\nu(x^\nu) \geq f(x). \tag{2}$$

The *Attouch-Wets (aw) distance* d , which quantifies the distance between usc functions in terms of a distance between their hypographs, metrizes hypo-convergence. Specifically, for $f, g \in \text{usc-fcns}(S)$, it is defined as

$$d(f, g) = \int_0^\infty d_\rho(f, g) e^{-\rho} d\rho,$$

where, for $\rho \geq 0$, the ρ -*aw-distance*

$$d_\rho(f, g) = \max_{z \in \rho B_\infty} \left| \text{dist}_\infty(z, \text{hypo } f) - \text{dist}_\infty(z, \text{hypo } g) \right|,$$

with $\text{dist}_\infty(z, A)$ being the usual point-to-set distance between a point $z \in S \times \mathbb{R}$ and a set $A \subset S \times \mathbb{R}$ under the sup-norm, $\rho B_\infty = B_\infty(0, \rho)$, with $B_\infty(\bar{z}, \rho) = \{z \in S \times \mathbb{R} \mid \|\bar{z} - z\|_\infty \leq \rho\}$ for any $\bar{z} \in S \times \mathbb{R}$. Since the meaning will be clear from the context, we also write $B_\infty(\bar{x}, \rho) = \{x \in S \mid \|\bar{x} - x\|_\infty \leq \rho\}$ with $\bar{x} \in S$. For any nonempty closed set $S \subset \mathbb{R}^n$, $(\text{usc-fcns}(S), d)$ is a complete separable metric space. Every closed and bounded subset $F \subset (\text{usc-fcns}(S), d)$ is compact. Moreover, for all $f, g \in \text{usc-fcns}(S)$,

$$\left| \text{dist}_\infty(0, \text{hypo } f) - \text{dist}_\infty(0, \text{hypo } g) \right| \leq d(f, g) \leq \max \left\{ \text{dist}_\infty(0, \text{hypo } f), \text{dist}_\infty(0, \text{hypo } g) \right\} + 1 \tag{3}$$

and, thus, if $f, g \geq 0$, then $0 \leq d(f, g) \leq 1$. We also see that a sufficient condition for F to be bounded is that there exists $(x, \alpha) \in S \times \mathbb{R}$ such that $f(x) \geq \alpha$ for all $f \in F$.

If not specified otherwise, the index ν runs over \mathbb{N} so that $x^\nu \rightarrow x$ means that the whole sequence $\{x^\nu, \nu \in \mathbb{N}\}$ converges to x . Let

$\mathcal{N}_\infty^\#$ be all the subsets of \mathbb{N} determined by subsequences,

i.e., $N \in \mathcal{N}_\infty^\#$ is an infinite collection of strictly increasing natural numbers. Thus, $\{x^\nu, \nu \in N\}$ is a subsequence of $\{x^\nu, \nu \in \mathbb{N}\}$; its convergence to x is noted by $x^\nu \xrightarrow{N} x$.

A similar development is available for functions defined on the metric space $(\text{usc-fcns}(S), \mathbf{d})$. However, we adopt a slightly different set-up that highlights the role of domains of definition. For $F, F^\nu \subset \text{usc-fcns}(S)$, the functions $\varphi^\nu : F^\nu \rightarrow \overline{\mathbb{R}}$ *epi-converge to* $\varphi : F \rightarrow \overline{\mathbb{R}}$ whenever

$$\begin{aligned} \forall N \in \mathcal{N}_\infty^\# \text{ and } f^\nu \in F^\nu \xrightarrow{N} f, \liminf_{\nu \in N} \varphi^\nu(f^\nu) \geq \varphi(f) \text{ if } f \in F \text{ and } \varphi^\nu(f^\nu) \xrightarrow{N} \infty \text{ otherwise} \\ \forall f \in F, \exists f^\nu \in F^\nu \rightarrow f \text{ with } \limsup \varphi^\nu(f^\nu) \leq \varphi(f). \end{aligned}$$

For $\varepsilon \geq 0$, ε - $\text{argmin}_{f \in F} \varphi(f) = \{f \in F \mid \varphi(f) \leq \inf_{g \in F} \varphi(g) + \varepsilon\}$, with the usual extended real-valued calculus in play when needed. We deviate slightly from the convention in [31] by setting ε - $\text{argmin}_{f \in F} \varphi(f) = F$ when $\varphi(f) = \infty$ for all $f \in F$. This is tenable because we restrain from extending functions to the whole space and ∞ is not assigned a special role in that regard. Consequently, we alleviate the need for checking that functions are finite at least somewhere, which in an infinite-dimensional setting may require excessively strong assumptions.

3 Piecewise Affine Approximations

In this section, we establish that every usc function on $S \subset \mathbb{R}^n$ can be approximated by piecewise affine functions with a particular structure under the additional assumption that S is convex. For $\rho \in [0, \infty)$ and $q \in \mathbb{N}$, let

$$\begin{aligned} \text{pa-fcns}_\rho^q(S) = \left\{ f : S \rightarrow [-\infty, \infty) \mid \exists a^k, b^k \in \mathbb{R}^n, \alpha_k, \beta_k \in \mathbb{R}, k = 1, \dots, q \text{ such that} \right. \\ \left. f(x) = \max_{k=1, \dots, q} [\langle a^k, x \rangle + \alpha_k] - \max_{k=1, \dots, q} [\langle b^k, x \rangle + \beta_k] \forall x \in S \cap \rho \mathbb{B}_\infty; f(x) = -\infty \text{ otherwise} \right\}. \end{aligned}$$

A function in $\text{pa-fcns}_\rho^q(S)$ is a difference of pointwise maxima of affine functions on $S \cap \rho \mathbb{B}_\infty$ and therefore is finite and continuous on that set. We say that $U \subset \text{usc-fcns}(S)$ is *hypo-dense* in $\text{usc-fcns}(S)$ if every $f \in \text{usc-fcns}(S)$ is the limit of a hypo-converging sequence $\{f^\nu \in U, \nu \in \mathbb{N}\}$.

3.1 Theorem (piecewise affine approximations). *Suppose that S is convex and $\rho^\nu \in [0, \infty)$ as well as $q^\nu \in \mathbb{N}$ tend to ∞ . Then,*

$$\bigcup_{\nu \in \mathbb{N}} \text{pa-fcns}_{\rho^\nu}^{q^\nu}(S) \text{ is hypo-dense in } \text{usc-fcns}(S).$$

Proof. Let $f \in \text{usc-fcns}(S)$. We construct a sequence in $\bigcup_{\nu \in \mathbb{N}} \text{pa-fcns}_{\rho^\nu}^{q^\nu}(S)$ that hypo-converges to f . Let $f^\nu : S \rightarrow \overline{\mathbb{R}}$ have $f^\nu(x) = \min\{f(x), \nu\}$ for all $x \in S$. Clearly, $f^\nu \rightarrow f$; recall that \rightarrow always denotes hypo-convergence when written between usc functions. For any $\nu \in \mathbb{N}$, f^ν is (upper) prox-bounded so that for every $\lambda > 0$, the (upper) Moreau-envelope $e_\lambda f^\nu : S \rightarrow \mathbb{R}$ of f^ν , which is given by

$$e_\lambda f^\nu(x) = \sup_{y \in S} f^\nu(y) - \frac{1}{2\lambda} \|y - x\|_2^2,$$

is finite and continuous. Moreover, $e_\lambda f^\nu \rightarrow f^\nu$ as $\lambda \searrow 0$ (see for example the discussion after Proposition 7.4 in [31]). Thus, there exists $\{\lambda^\nu > 0, \nu \in \mathbb{N}\} \rightarrow 0$ such that

$$e_{\lambda^\nu} f^\nu \rightarrow f.$$

Next, we define $\varphi^\nu : S \rightarrow \overline{\mathbb{R}}$ as

$$\varphi^\nu(x) = e_{\lambda^\nu} f^\nu(x) \quad \forall x \in S \cap \rho^\nu \mathcal{B}_\infty \text{ and } \varphi^\nu(x) = -\infty \text{ otherwise.}$$

Since $\rho^\nu \mathcal{B}_\infty \rightarrow \mathbb{R}^n$, we also have that $\varphi^\nu \rightarrow f$.

Every real-valued continuous function on a convex compact subset of \mathbb{R}^n is the limit in the sup-norm of finite-valued dc functions defined on the same set; see for example [21, Prop. 2.3]. Consequently, for every $\nu \in \mathbb{N}$, there exist convex functions $\{g_\mu^\nu, h_\mu^\nu : S \rightarrow \overline{\mathbb{R}}, \mu \in \mathbb{N}\}$, finite on $S \cap \rho^\nu \mathcal{B}_\infty$, with the property that

$$\sup_{x \in S \cap \rho^\nu \mathcal{B}_\infty} |\varphi^\nu(x) - [g_\mu^\nu(x) - h_\mu^\nu(x)]| \rightarrow 0 \text{ as } \mu \rightarrow \infty.$$

Let $\psi_\mu^\nu : S \rightarrow \overline{\mathbb{R}}$ be defined by

$$\psi_\mu^\nu(x) = g_\mu^\nu(x) - h_\mu^\nu(x) \text{ for } x \in S \cap \rho^\nu \mathcal{B}_\infty \text{ and } \psi_\mu^\nu(x) = -\infty \text{ otherwise.}$$

Since already $\varphi^\nu \rightarrow f$, we can construct $\{\mu^\nu \in \mathbb{N}, \nu \in \mathbb{N}\} \rightarrow \infty$ as $\nu \rightarrow \infty$ such that $\psi_{\mu^\nu}^\nu \rightarrow f$. Let $\psi^\nu = \psi_{\mu^\nu}^\nu$.

The convex functions $\{g_{\mu^\nu}^\nu, h_{\mu^\nu}^\nu, \nu \in \mathbb{N}\}$ are lsc and proper. Let $g^\nu = g_{\mu^\nu}^\nu$ and $h^\nu = h_{\mu^\nu}^\nu$. Consequently, for every $\nu \in \mathbb{N}$,

$$g^\nu(x) = \sup_{(a,\alpha) \in A(g^\nu)} \{\langle a, x \rangle + \alpha\} \text{ and } h^\nu(x) = \sup_{(a,\alpha) \in A(h^\nu)} \{\langle a, x \rangle + \alpha\} \text{ for } x \in S,$$

where for $u : S \rightarrow \overline{\mathbb{R}}$,

$$A(u) = \{(a, \alpha) \in \mathbb{R}^n \times \mathbb{R} \mid \langle a, x \rangle + \alpha \leq u(x) \quad \forall x \in S\}.$$

Since $A(g^\nu), A(h^\nu) \subset \mathbb{R}^{n+1}$, which is separable, there exist increasing sets $\{A^\mu(g^\nu), A^\mu(h^\nu), \mu \in \mathbb{N}\}$, each with finite cardinality, such that

$$\bigcup_{\mu \in \mathbb{N}} A^\mu(g^\nu) \text{ is dense in } A(g^\nu) \text{ and } \bigcup_{\mu \in \mathbb{N}} A^\mu(h^\nu) \text{ is dense in } A(h^\nu).$$

For $\nu, \mu \in \mathbb{N}$, we define $\tilde{g}_\mu^\nu, \tilde{h}_\mu^\nu : S \rightarrow \overline{\mathbb{R}}$ by setting

$$\tilde{g}_\mu^\nu(x) = \max_{(a,\alpha) \in A^\mu(g^\nu)} \{\langle a, x \rangle + \alpha\} \text{ and } \tilde{h}_\mu^\nu(x) = \max_{(a,\alpha) \in A^\mu(h^\nu)} \{\langle a, x \rangle + \alpha\} \text{ for } x \in S \cap \rho^\nu \mathcal{B}_\infty,$$

and $\tilde{g}_\mu^\nu(x) = \tilde{h}_\mu^\nu(x) = \infty$ otherwise. The characterization of hypo-convergence in (1) and (2) enables us to conclude that for all $\nu \in \mathbb{N}$,

$$-\tilde{g}_\mu^\nu \rightarrow -g^\nu \text{ and } -\tilde{h}_\mu^\nu \rightarrow -h^\nu \text{ as } \mu \rightarrow \infty,$$

with pointwise convergence holding as well on $S \cap \rho^\nu \mathcal{B}_\infty$.

Let $\tilde{\psi}_\mu^\nu : S \rightarrow \overline{\mathbb{R}}$ be defined by

$$\tilde{\psi}_\mu^\nu(x) = \tilde{g}_\mu^\nu(x) - \tilde{h}_\mu^\nu(x) \text{ for } x \in S \cap \rho^\nu \mathcal{B}_\infty \text{ and } \tilde{\psi}_\mu^\nu(x) = -\infty \text{ otherwise.}$$

If $\tilde{\psi}_\mu^\nu \rightarrow \psi^\nu$ as $\mu \rightarrow \infty$, then we can construct $\{\mu^\nu \in \mathbb{N}, \nu \in \mathbb{N}\} \rightarrow \infty$ such that

$$\tilde{\psi}_{\mu^\nu}^{\nu} \rightarrow f$$

because already $\psi^\nu \rightarrow f$. Since $\tilde{\psi}_{\mu^\nu}^{\nu} \in \text{pa-fcns}_{\rho^\nu}^{q^\nu}(S)$, with q^ν being the largest cardinality of $A^{\mu^\nu}(g^\nu)$ and of $A^{\mu^\nu}(h^\nu)$, the conclusion will follow.

It only remains to establish that $\tilde{\psi}_\mu^\nu \rightarrow \psi^\nu$ as $\mu \rightarrow \infty$. For this purpose, we again leverage the characterization of hypo-convergence in (1) and (2). Suppose that $x^\mu \in S \cap \rho^\nu \mathcal{B}_\infty \rightarrow x$. Then,

$$\begin{aligned} \limsup_\mu (\tilde{g}_\mu^\nu(x^\mu) - \tilde{h}_\mu^\nu(x^\mu)) &= \limsup_\mu \tilde{g}_\mu^\nu(x^\mu) - \liminf_\mu \tilde{h}_\mu^\nu(x^\mu) \\ &\leq \limsup_\mu g^\nu(x^\mu) - h^\nu(x) \leq g^\nu(x) - h^\nu(x), \end{aligned}$$

where we use the facts that \tilde{g}_μ^ν lower bounds g^ν , $-\tilde{h}_\mu^\nu \rightarrow -h^\nu$, and g^ν is continuous on $S \cap \rho^\nu \mathcal{B}_\infty$. Also, for $x \in S \cap \rho^\nu \mathcal{B}_\infty$,

$$\liminf_\mu (\tilde{g}_\mu^\nu(x) - \tilde{h}_\mu^\nu(x)) = \liminf_\mu \tilde{g}_\mu^\nu(x) - \limsup_\mu \tilde{h}_\mu^\nu(x) \geq g^\nu(x) - h^\nu(x).$$

These inequalities are trivially satisfied for sequences outside of $S \cap \rho^\nu \mathcal{B}_\infty$. Hence, the assertion is established. \square

We illustrate the usefulness of the theorem in the solution of *function identification problems* of the form

$$\text{(FIP)} \quad \min_{f \in F} \varphi(f), \text{ where } F \subset \text{usc-fcns}(S) \text{ and } \varphi : F \rightarrow \overline{\mathbb{R}}.$$

These problems arise in nonparametric estimation, spatial statistics, and curve fitting; see [34] for applications to estimation of financial curves, electricity demand, commodity prices, and uncertainty in physical systems. For example, if F is a class of n -dimensional probability density functions and $\varphi(f) = -\frac{1}{m} \sum_{j=1}^m \log f(x^j)$, then any minimizer of (FIP) is a maximum likelihood estimate based on the data $x^1, \dots, x^m \in S$. When $\varphi(f) = \frac{1}{m} \sum_{j=1}^m (y^j - f(x^j))^2$, a minimizer furnishes a least-squares fit of the data $\{(x^j, y^j) \in S \times \mathbb{R}, j = 1, \dots, m\}$ over the class F . We refer to [35, 38] for numerous examples. Further unexplored applications for usc functions and their piecewise affine approximations may arise in stochastic and robust optimization where problems can be formulated over spaces of decision rules and be approximated using polynomial and piecewise polynomial functions [2], finite collections of policies [17], and linear decision rules, possibly in a higher dimensional spaces [12]. In special cases of (FIP), such as when F consists of concave functions only, one might be able to reformulate the problem as an equivalent finite-dimensional one; see [9] for the context of maximum likelihood estimation over the log-concave class and [41] for least-squares regression over convex functions. However, this is not possible in general and we need to settle for approximations.

For $\rho^\nu \in [0, \infty)$ and $q^\nu \in \mathbb{N}$, we consider the *approximating function identification problem*

$$(\text{FIP})^\nu \quad \min_{f \in F^\nu} \varphi(f), \text{ where } F^\nu = F \cap \text{pa-fcns}_{\rho^\nu}^{q^\nu}(S).$$

Every function in $\text{pa-fcns}_{\rho^\nu}^{q^\nu}(S)$ is described by $2q^\nu(n+1)$ parameters¹. Thus, $(\text{FIP})^\nu$ is equivalent to a finite-dimensional optimization problem with the same number of variables. The number grows only linearly in n , which makes the approach promising for high-dimensional problems. In comparison, approximations based on epi-splines (see [35, 38]) require a preselected partition of S not easily decided on in a computationally tractable manner beyond four or five dimensions. The piecewise approximations in $(\text{FIP})^\nu$ are mesh free, with the domain of each affine component adapting to the problem at hand. Thus, we expect to be able to identify and leverage low-dimensional structures, if present, when solving (FIP) by means of $(\text{FIP})^\nu$.

It is apparent that an application may demand approximations also of the objective function in (FIP) , which we address in Section 5 for the central case of stochastic optimization where $\varphi = \mathbb{E}[\psi(\boldsymbol{\xi}, \cdot)]$; the expectation with respect to the distribution of a random element $\boldsymbol{\xi}$ is denoted by \mathbb{E} . Here, we concentrate on the application of Theorem 3.1 to justify $(\text{FIP})^\nu$.

Let $\varphi^\nu : F^\nu \rightarrow \overline{\mathbb{R}}$ be the function defined by $\varphi^\nu(f) = \varphi(f)$ for $f \in F^\nu$. Suppose that S is convex, F is a nonempty and *solid* subset of $(\text{usc-fcns}(S), \mathcal{d})$, i.e., $F = \text{cl}(\text{int } F)$, and $\varphi : F \rightarrow \overline{\mathbb{R}}$ is continuous on F . Then, a standard argument (see for example the proof of Theorem 3.16 in [35]) in conjunction with Theorem 3.1 establishes that

$$\varphi^\nu \text{ epi-converges to } \varphi \text{ provided that } \rho^\nu, q^\nu \rightarrow \infty.$$

Thus, when $\{\varepsilon^\nu \geq 0, \nu \in \mathbb{N}\} \rightarrow 0$,

$$\text{OutLim} \left(\varepsilon^\nu\text{-argmin}_{f \in F^\nu} \varphi(f) \right) \subset \text{argmin}_{f \in F} \varphi(f),$$

which can be deduced, for example, from [33]. The constraint qualification that F is solid cannot be relaxed without introducing some other assumption. Obviously, $F \cap \text{pa-fcns}_{\rho^\nu}^{q^\nu}(S)$ can, in general, be empty for all ν , but when F is solid this is ruled out.

In view of this discussion, the challenge of solving an infinite-dimensional function identification problem from the broad class (FIP) is shifted to that of obtaining a near-minimizer of a finite-dimensional problem. Of course, the difficulty of that task depends on the specific properties of φ and F . Typically, $(\text{FIP})^\nu$ would be nonconvex, but the special affine structure of functions in $\text{pa-fcns}_{\rho^\nu}^{q^\nu}(S)$ is bound to be important in developing computational procedures. Initial efforts in that direction are already found in [8], which presents several algorithms with guarantees to obtain at least certain stationary points and numerical results from nonparametric least-squares regression, as well as in [28], which approximates functions in up to $n = 41$ dimensions using piecewise affine functions in the context of nonparametric support vector machines. The nonconvexity of $(\text{FIP})^\nu$ encountered in [28] appears to be only moderately challenging and handled by common randomization strategies.

¹We stress that ν is an index and not the power of q .

4 Covering Numbers

It is well known that every bounded $F \subset (\text{usc-fcns}(S), \mathcal{d})$ has a finite cover by virtue of being totally bounded. However, this is not sufficient to establish certain rates of convergence results for sample average approximations of stochastic optimization problems of the form $\min_{f \in F} \mathbb{E}[\psi(\boldsymbol{\xi}, f)]$. It is usually necessary to bound for all $\varepsilon > 0$ the *covering number* of F , denoted by $N(F, \varepsilon)$, which is the smallest number of closed \mathcal{d} -balls of radius ε needed to cover F . We next provide such a bound and show that it is nearly sharp. Section 5 applies the result to establish rates of convergence of minimizers of stochastic optimization problems.

We start by recording a useful estimate of the hypo-distance. For $f, g \in \text{usc-fcns}(S)$ and $\rho \geq 0$, we define the auxiliary quantity

$$\hat{\mathcal{d}}_\rho(f, g) = \inf \left\{ \tau \geq 0 \mid \begin{array}{l} \sup_{y \in \mathcal{B}_\infty(x, \tau)} g(y) \geq \min\{f(x), \rho\} + \tau, \forall x \in \rho \mathcal{B}_\infty \text{ with } f(x) \geq -\rho \\ \sup_{y \in \mathcal{B}_\infty(x, \tau)} f(y) \geq \min\{g(x), \rho\} + \tau, \forall x \in \rho \mathcal{B}_\infty \text{ with } g(x) \geq -\rho \end{array} \right\}.$$

As the notation indicates, $\hat{\mathcal{d}}_\rho$ is closely related to \mathcal{d}_ρ (see Proposition 3.1 in [33]) and therefore also to \mathcal{d} . We record the relevant properties next.

4.1 Lemma For $f, g \in \text{usc-fcns}(S)$ and $\rho \geq 0$,

$$e^{-\rho} \hat{\mathcal{d}}_\rho(f, g) \leq \mathcal{d}(f, g) \leq (1 - e^{-\rho}) \hat{\mathcal{d}}_{2\rho + \delta}(f, g) + e^{-\rho}(\delta + \rho + 1),$$

where $\delta = \max\{\text{dist}_\infty(0, \text{hypo } f), \text{dist}_\infty(0, \text{hypo } g)\}$.

Proof. The results can be deduced from Propositions 3.1 and 3.2 in [33]. □

As mentioned in Section 1, epi-splines [35, 33] furnish a dense subset of classes of semicontinuous functions and associated error bounds are known. We leverage these results here. Although the piecewise affine functions of Section 3 are also dense in the usc functions, they have unknown error and cannot presently serve as the basis for the construction in the proof of the next theorem. This is anyhow less critical as we see through a lower bound result (Theorem 4.4 below) that the obtained upper bound on covering numbers is within a logarithmic factor of being sharp.

For any $f : S \rightarrow \overline{\mathbb{R}}$ and $x \in S$, let $\liminf_{\bar{x} \rightarrow x} f(\bar{x}) = \lim_{\delta \downarrow 0} \inf_{\bar{x} \in \mathcal{B}_\infty(x, \delta)} f(\bar{x})$. Epi-splines are defined in terms a finite collection of subsets of S . A finite collection R_1, R_2, \dots, R_K of open subsets² of S is a *partition* of S if $\cup_{k=1}^K \text{cl } R_k = S$ and $R_k \cap R_l = \emptyset$ for all $k \neq l$. Specifically, an *epi-spline* $s : S \rightarrow \mathbb{R}$, with partition $\mathcal{R} = \{R_1, \dots, R_K\}$ of S , is a function that

$$\begin{array}{l} \text{on each } R_k, k = 1, \dots, K, \text{ takes a constant real number as value,} \\ \text{and for every } x \in S, \text{ has } s(x) = \liminf_{x' \rightarrow x} s(x'). \end{array}$$

The family of all such epi-splines is denoted by $\text{e-spl}(\mathcal{R})$. Epi-splines are lsc by construction and approximate lsc functions in the sense of epi-convergence. Since the present setting involves usc functions

²Recall that “open” here is according to the metric space $(S, \|\cdot - \cdot\|_\infty)$.

and hypo-convergence, we “reorientation” and introduce minus in some expressions. We refer to [35, 33] for further information and extensions that go beyond these zeroth order epi-splines and also beyond \mathbb{R}^n .

4.2 Proposition *For a partition $\mathcal{R} = \{R_1, \dots, R_K\}$ of S and $\rho \geq 0$, we have that for every $f \in \text{usc-fcns}(S)$, there exists an $s \in \text{e-spl}(\mathcal{R})$ such that*

$$\hat{d}_\rho(f, -s) \leq \mu_\rho(\mathcal{R}) = \inf \{ \tau \geq 0 \mid R_k \subset \mathbb{B}_\infty(x, \tau) \text{ for all } x \in \rho \mathbb{B}_\infty \text{ and } k \text{ satisfying } x \in \text{cl } R_k \}.$$

If $\mu_\rho(\mathcal{R}) \leq \rho$, then s can be taken to satisfy $-\rho' \leq s(x) \leq \max\{-\rho', \min[\rho', -f(x)]\}$ for any $\rho' > \rho$ and $x \in S$.

Proof. The first part of the proposition is a direct application of [33, Theorem 5.9]. The fact that s can be taken to satisfy $-\rho' \leq s(x) \leq \max\{-\rho', \min[\rho', -f(x)]\}$ for any $\rho' > \rho$ follows from an examination of that theorem’s proof. \square

The quantity $\mu_\rho(\mathcal{R})$ is the *meshsize* of $\mathcal{R} = \{R_1, \dots, R_K\}$ and, essentially, quantifies the size of the largest R_k .

4.3 Theorem (covering numbers). *For every bounded subset F of $(\text{usc-fcns}(S), \mathbf{d})$, there exist $c \geq 0$ and $\bar{\varepsilon} > 0$ (both independent of n , the dimension of S) such that*

$$\log N(F, \varepsilon) \leq \left(\frac{c}{\varepsilon}\right)^n \left(\log \frac{1}{\varepsilon}\right)^{n+1} \text{ for all } \varepsilon \in (0, \bar{\varepsilon}].$$

Proof. Since F is bounded, there exists an $r > 0$ such that $\text{dist}_\infty(0, \text{hypo } f) \leq r$ for all $f \in F$. Let $\gamma_1, \gamma_2, \gamma_3 > 0$ be such that $\gamma_1 + \gamma_2 + \gamma_3 = 1$. Set $\bar{\varepsilon} \in (0, 1)$ such that

$$\frac{2(r+1)}{r} \left[\log \frac{1}{\varepsilon} + \log \frac{1}{\gamma_1} + \frac{r}{2} + \log(r+1) \right] - 1 > \gamma_2 \varepsilon \text{ for all } \varepsilon \in (0, \bar{\varepsilon}]. \quad (4)$$

Fix $\varepsilon \in (0, \bar{\varepsilon}]$ and define ρ to be the expression on the left-hand side of (4). We next construct a partition of S and set $\omega > 1$ and

$$\nu = \left\lceil \frac{2\omega\rho}{\gamma_2\varepsilon} \right\rceil,$$

where $\lceil a \rceil$ is the smallest integer no smaller than a . The partition is obtained by dividing the box $[-\omega\rho, \omega\rho]^n \subset \mathbb{R}^n$ into ν^n boxes of equal size and then intersecting with S . Let $K = \nu^n + 1$. Specifically, for $k = 1, 2, \dots, \nu^n$, set

$$R_k = \text{int} \left(S \cap \prod_{i=1}^n (l_i^k, u_i^k) \right), \text{ with } l_i^k = 2(k-1)\omega\rho/\nu - \omega\rho \text{ and } u_i^k = l_i^k + 2\omega\rho/\nu$$

so that $\cup_{k=1}^{K-1} \text{cl } R_k = S \cap [-\omega\rho, \omega\rho]^n$. Also, $R_K = \text{int}(S \setminus [-\omega\rho, \omega\rho]^n)$. Again, we recall that the interior and closure are taken relative to $(S, \|\cdot\|_\infty)$. We denote by $\mathcal{R} = \{R_1, \dots, R_K\}$ this partition. Clearly, $\mu_\rho(\mathcal{R}) = 2\omega\rho/\nu$. Next, we consider a discretization of parts of the range of functions and set

$$m = \left\lceil \frac{\omega\rho}{\gamma_3\varepsilon} \right\rceil + 1.$$

The points $\sigma_j = -\omega\rho + 2(j-1)\omega\rho/(m-1)$, $j = 1, 2, \dots, m$, discretize the interval $[-\omega\rho, \omega\rho]$. Let F_0 be the collection of piecewise constant functions on \mathcal{R} with values in $\{\sigma_1, \dots, \sigma_m\}$ defined as follows. If $f \in F_0$, then for every $k \in \{1, \dots, K\}$ there exists a $j_k \in \{1, \dots, m\}$ such that $f(x) = \sigma_{j_k}$ for $x \in R_k$ and $f(x) = \lim_{\delta \downarrow 0} \sup_{y \in B_\infty(x, \delta)} f(y)$ otherwise. By construction, f is usc. Obviously, F_0 contains m^K functions. We now show that every $f \in F$ has $d(f, f_0) \leq \varepsilon$ for some $f_0 \in F_0$.

Let $f \in F$ be arbitrary. By Proposition 4.2 and the fact that $\mu_\rho(\mathcal{R}) = 2\omega\rho/\nu \leq \gamma_2\varepsilon < \rho$, there exists $s \in \text{e-spl}(\mathcal{R})$ such that

$$\hat{d}_\rho(f, -s) \leq \mu_\rho(\mathcal{R}) \text{ and } -\omega\rho \leq s(x) \leq \max\{-\omega\rho, \min[\omega\rho, -f(x)]\} \text{ for } x \in S.$$

Since $\text{dist}(0, \text{hypo } f) \leq r$, there exists $x \in rB_\infty$ such that $f(x) \geq -r$. Consequently, $-s(x) \geq \min\{\omega\rho, \max[-\omega\rho, f(x)]\} \geq -r$. So we also have that $\text{dist}_\infty(0, \text{hypo } -s) \leq r$.

Since $\varepsilon, \gamma_1 \leq 1$,

$$\rho \geq \frac{2(r+1)}{r} \left[\frac{r}{2} + \log(r+1) \right] - 1 = r + \frac{2(r+1)}{r} \log(r+1) \geq r.$$

Thus, using the notation $\bar{\rho} = (\rho - r)/2$, Lemma 4.1 gives that

$$d(f, -s) \leq \hat{d}_\rho(f, -s) + e^{-\bar{\rho}}(r + \bar{\rho} + 1) \leq \mu_\rho(\mathcal{R}) + e^{-\bar{\rho}}(r + \bar{\rho} + 1) = 2\omega\rho/\nu + e^{-\bar{\rho}}(r + \bar{\rho} + 1).$$

In view of [35, Theorem 3.17], there exists $f_0 \in F_0$ such that $d(-s, f_0) \leq \omega\rho/(m-1)$ since we can select f_0 such that $|s(x) - f_0(x)| \leq \omega\rho/(m-1)$ for all $x \in S$. The triangle inequality then gives that

$$d(f, f_0) \leq \omega\rho/(m-1) + 2\omega\rho/\nu + e^{-\bar{\rho}}(r + \bar{\rho} + 1). \quad (5)$$

It remains to show that the right-hand side is no greater than ε . We start with the last term in (5). By concavity of the log-function, we have that

$$\log\left(\frac{1}{2}(\rho + r) + 1\right) \leq \log(r+1) + \frac{\rho - r}{2r + 2}.$$

Consequently,

$$\begin{aligned} \log[e^{-\bar{\rho}}(r + \bar{\rho} + 1)] &= \frac{1}{2}(r - \rho) + \log\left(\frac{1}{2}(\rho + r) + 1\right) \leq \frac{1}{2}(r - \rho) + \log(r+1) + \frac{\rho - r}{2r + 2} \\ &= \frac{r}{2} - \frac{r(\rho + 1)}{2(r+1)} + \log(r+1) = \log \gamma_1 \varepsilon, \end{aligned}$$

where the last equality follows from inserting the expression for ρ . Thus, $e^{-\bar{\rho}}(r + \bar{\rho} + 1) \leq \gamma_1 \varepsilon$. We next examine the second term on the right-hand side of (5). Inserting the expression for ν , we obtain that

$$\frac{2\omega\rho}{\nu} \leq \gamma_2 \varepsilon.$$

Finally, we consider the first term on the right-hand side of (5). In view of the definition of m , we have that

$$\frac{\omega\rho}{m-1} \leq \gamma_3 \varepsilon.$$

Thus, $d(f, f_0) \leq \varepsilon$ and we have established that d -balls with radius ε and centered at points in F_0 cover F . The logarithm of the number of functions in F_0 is $(\nu^n + 1) \log m$. At this point, the order of the result is immediate. A possible expression for the constant c is obtained as follows. Let $c_1 = 2(r + 1)/r$ and

$$c_2 = \frac{2(r + 1)}{r} \left[\log \frac{1}{\gamma_1} + \frac{r}{2} + \log(r + 1) \right] - 1.$$

Thus, $\rho = c_1 \log \varepsilon^{-1} + c_2$. Moreover, let $c_3 = 2\omega/\gamma_2$ and $c_4 = \omega/\gamma_3$. Using these expressions, we find that

$$(\nu^n + 1) \log m \leq \left[\left(c_1 c_3 + \frac{c_2 c_3 + 1}{\log \bar{\varepsilon}^{-1}} \right)^n \left(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon} \right)^n + 1 \right] \log \left[\left(c_1 c_4 + \frac{c_2 c_4 + 2}{\log \bar{\varepsilon}^{-1}} \right) \frac{1}{\varepsilon} \log \frac{1}{\varepsilon} \right].$$

Let

$$c_5 = c_1 c_3 + \frac{c_2 c_3 + 1}{\log \bar{\varepsilon}^{-1}} \text{ and } c_6 = c_1 c_4 + \frac{c_2 c_4 + 2}{\log \bar{\varepsilon}^{-1}}.$$

We then find that

$$(\nu^n + 1) \log m \leq c_7^n \left(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon} \right)^n \left[\log c_6 + \log \frac{1}{\varepsilon} + \log \log \frac{1}{\varepsilon} \right], \text{ where } c_7 = c_5 + \frac{1}{\bar{\varepsilon}^{-1} \log \bar{\varepsilon}^{-1}}.$$

Using the fact that $\log \log \varepsilon^{-1} / \log \varepsilon^{-1} \leq e^{-1}$ for $\varepsilon \in (0, 1)$, we obtain

$$(\nu^n + 1) \log m \leq c_7^n \left[\frac{\log c_6}{\log \bar{\varepsilon}^{-1}} + 1 + e^{-1} \right] \frac{1}{\varepsilon^n} \left(\log \frac{1}{\varepsilon} \right)^{n+1}, \quad (6)$$

which gives a particular expression for c in the theorem statement. Since the choice of $\bar{\varepsilon}$ is independent of n , this c is independent of n . For example, for $\bar{\varepsilon} = 0.01$, $\omega = 0.00000001$, $r = 3.22$, then $c_7 = 13.5$ and the term in brackets in (6) evaluates to 2.3. \square

Although a comparison to the classical result of $O(\varepsilon^{-n})$ for Lipschitz continuous functions on bounded subsets, which goes back to [23] (see for example [44, Theorem 2.7.1]), is not entirely relevant due to the different settings, we note that our bound is only slightly worse (a logarithmic term) for larger families of usc functions. We do not require any bound on the variation of the functions and allow functions defined on all of \mathbb{R}^n , possibly extended real-valued. Still, the entropy integral³ $\int_0^{\bar{\varepsilon}} \sqrt{\log N(F, \varepsilon)} d\varepsilon$ is finite only for $n = 1$ and, therefore, these families are “large,” and increasingly so as n grows.

4.4 Theorem (covering numbers; lower bound). *For every $n \in \mathbb{N}$, there exist a bounded subset $F \subset \text{usc-fns}(\mathbb{R}^n)$ and corresponding $c \geq 0$ and $\bar{\varepsilon} > 0$ (independent of n) such that*

$$\log N(F, \varepsilon) \geq \left(\frac{c}{\varepsilon} \right)^n \log \frac{1}{\varepsilon} \text{ for all } \varepsilon \in (0, \bar{\varepsilon}].$$

³For the significance of entropy integrals we refer to [44].

Proof. See the appendix. □

In comparison with the upper bound of Theorem 4.3, we see that the lower bound differs by a logarithmic factor only and therefore the upper bound is nearly sharp. The size of the bounded set F in Theorem 4.4 does not have to be large. In fact, an examination of the proof reveals that F might be selected to have $d(0, f) \leq r$ for all $f \in F$, with r being only slightly above one. Here, 0 is the function in $\text{usc-fcns}(\mathbb{R}^n)$ that is identical to zero.

The proof of Theorem 4.4 constructs a collection of functions which is finite on a grid of points in $[0, \rho]^n$, with $\rho > 0$ and grid points spaced roughly ε apart. At each of these grid points, a function takes on a value among a set of discretized values between $-\rho$ and 0, again spaced roughly ε apart. Outside these grid points, the functions are assigned $-\infty$. It is clear that the number of such functions is $(\rho/\varepsilon)^\nu$, where $\nu = (\rho/\varepsilon)^n$. Thus, its logarithm is of the order $O(\varepsilon^{-n} \log \varepsilon^{-1})$.

5 Applications to Stochastic Optimization and Statistical Estimation

Suppose that $(\Xi, \mathcal{A}, \mathbb{P})$ is a complete probability space, $F \subset \text{usc-fcns}(S)$ is closed, and $\psi : \Xi \times F \rightarrow \overline{\mathbb{R}}$ is a function with suitable properties as discussed below. We denote by boldface, for example $\boldsymbol{\xi}$, random elements with values in Ξ . Then, a *function identification problem under uncertainty* takes the form

$$\text{(FIP-U)} \quad \min_{f \in F} \mathbb{E}[\psi(\boldsymbol{\xi}, f)] = \int \psi(\boldsymbol{\xi}, f) d\mathbb{P}(\boldsymbol{\xi}).$$

Section 3 furnishes some instances of ψ in the context of probability density estimation and regression, see also below, but there are numerous other examples.

A *sample average approximation* of the problem leverages a sample $\boldsymbol{\xi}^1, \boldsymbol{\xi}^2, \dots$ of independent random elements, each with values in Ξ and distributed according to \mathbb{P} , and leads to the *approximating problem*

$$\text{(FIP-U)}^\nu \quad \min_{f \in F} \frac{1}{\nu} \sum_{j=1}^{\nu} \psi(\boldsymbol{\xi}^j, f).$$

Under mild assumptions (see Proposition 5.1 below), minimizers of $(\text{FIP-U})^\nu$ tend to those of (FIP-U) almost surely. However, a main challenge in the justification of such an approach is to quantify the *rate* with which the error in solutions of $(\text{FIP-U})^\nu$ vanishes as ν grows. Before stating the results that rely on the covering numbers of the previous section, we formalize the setting. The following definitions and facts are well known; see for example⁴ [31, Ch. 14].

We say that $\psi : \Xi \times F \rightarrow \overline{\mathbb{R}}$ is a *random lsc function* if for all $\xi \in \Xi$, $\psi(\xi, \cdot)$ is lsc as a function on the metric space (F, d) and ψ is measurable with respect to the product sigma-algebra⁵ on $\Xi \times F$. A random lsc function $\psi : \Xi \times F \rightarrow \overline{\mathbb{R}}$ is *locally inf-integrable* if⁶

$$\forall f \in F \exists \rho > 0 \text{ such that } \int \inf_{g \in F} \{\psi(\xi, g) \mid d(g, f) \leq \rho\} d\mathbb{P}(\xi) > -\infty.$$

⁴This reference states results only for finite dimensions, but since (F, d) is a complete separable metric space, with compact balls, the proofs of the required results carry over nearly verbatim.

⁵On (F, d) , we adopt the Borel sigma-algebra.

⁶For measurable $h : \Xi \rightarrow \overline{\mathbb{R}}$, $\int h(\xi) d\mathbb{P}(\xi) = \int \max\{0, h(\xi)\} d\mathbb{P}(\xi) - \int \max\{0, -h(\xi)\} d\mathbb{P}(\xi)$, with $\infty - \infty = \infty$.

If $\psi : \Xi \times F \rightarrow \overline{\mathbb{R}}$ is a locally inf-integrable random lsc function, then $f \mapsto \mathbb{E}[\psi(\boldsymbol{\xi}, f)]$ is well-defined, always greater than $-\infty$, and lsc.

5.1 Proposition *Suppose that $(\Xi, \mathcal{A}, \mathbb{P})$ is a complete probability space, $F \subset \text{usc-fcns}(S)$ is closed, and $\psi : \Xi \times F \rightarrow \overline{\mathbb{R}}$ is an inf-integrable random lsc function. If $\boldsymbol{\xi}^1, \boldsymbol{\xi}^2, \dots$ is a sequence of independent random elements each with values in Ξ and distribution \mathbb{P} and $\{\varepsilon^\nu \geq 0, \nu \in \mathbb{N}\} \rightarrow 0$, then, almost surely,*

$$\text{OutLim} \left(\varepsilon^\nu\text{-argmin}_{f \in F} \frac{1}{\nu} \sum_{j=1}^{\nu} \psi(\boldsymbol{\xi}^j, f) \right) \subset \text{argmin}_{f \in F} \mathbb{E}[\psi(\boldsymbol{\xi}^1, f)].$$

Moreover, if F is bounded, then almost surely

$$\lim \left(\inf_{f \in F} \frac{1}{\nu} \sum_{j=1}^{\nu} \psi(\boldsymbol{\xi}^j, f) \right) = \inf_{f \in F} \mathbb{E}[\psi(\boldsymbol{\xi}^1, f)] > -\infty.$$

Proof. This is a consequence of a law of large numbers for lsc functions and epi-convergence; see for example Proposition 7.1 in [38]. \square

In the following, we assume that $F \subset \text{usc-fcns}(S)$ is closed and bounded as it results in some simplifications. In particular, (F, \mathcal{d}) then becomes a *compact* metric space. The assumption is anyhow minor as it is often acceptable in applications to impose a lower bound on the functions in $\text{usc-fcns}(S)$ under considerations; see the remark in conjunction with (3).

The *excess* of a set $F_1 \subset F$ over a set $F_2 \subset F$ is given by

$$\text{exs}(F_1, F_2) = \sup_{f \in F_1} \text{dist}(f, F_2) \text{ if } F_1, F_2 \text{ are nonempty,}$$

$\text{exs}(F_1, F_2) = \infty$ if F_1 nonempty and F_2 empty, and $\text{exs}(F_1, F_2) = 0$ otherwise. Here, $\text{dist}(f, F_2) = \inf_{g \in F_2} \mathcal{d}(f, g)$ is the usual point-to-set distance in (F, \mathcal{d}) . The Pompeiu-Hausdorff distance $\mathbb{H}(F_1, F_2) = \max\{\text{exs}(F_1, F_2), \text{exs}(F_2, F_1)\}$. Let the *level sets* of any $\varphi : F \rightarrow \overline{\mathbb{R}}$ be denoted by

$$\text{lev}_{\leq \delta} \varphi = \{f \in F \mid \varphi(f) \leq \delta\}.$$

If $\psi : \Xi \times F \rightarrow \overline{\mathbb{R}}$ is a random lsc function and $F_1, F_2 \subset F$ are closed, then

$$\xi \mapsto \text{exs} \left(\varepsilon\text{-argmin}_{f \in F_1} \psi(\xi, f), F_2 \right) \text{ and } \xi \mapsto \text{exs} \left(F_2, \text{lev}_{\leq \delta} \psi(\xi, \cdot) \right)$$

are random variables on $(\Xi, \mathcal{A}, \mathbb{P})$ for any $\varepsilon \geq 0$ and $\delta \in \mathbb{R}$.

When considering sample average approximations with sample size ν , the relevant probability space is the ν -fold product space formed by $(\Xi, \mathcal{A}, \mathbb{P})$; the above definitions apply also for this probability space. The *sample average function*

$$((\xi^1, \dots, \xi^\nu), f) \mapsto \frac{1}{\nu} \sum_{j=1}^{\nu} \psi(\xi^j, f)$$

is then a random lsc function on the product probability space provided that $\psi > -\infty$ is a random lsc function on $(\Xi, \mathcal{A}, \mathbb{P})$. Since this is the case below, the following probabilistic statements are meaningful. The measure on the product probability space is denoted by P^ν and the sample space by Ξ^ν .

We also need a quantitative result about differences between minimizers and related quantities. The following result improves on [33, Thms. 4.3 and 4.5]; see also [10] for related results in the convex setting. We denote by $\mathcal{B}(f, \rho) = \{g \in \text{usc-fcns}(S) \mid \mathcal{d}(f, g) \leq \rho\}$.

5.2 Proposition *For a closed and bounded $F \subset \text{usc-fcns}(S)$, let $F_1, F_2 \subset F$ be nonempty and $\varphi_1 : F_1 \rightarrow (-\infty, \infty]$ as well as $\varphi_2 : F_2 \rightarrow (-\infty, \infty]$ be lsc functions on the metric space (F, \mathcal{d}) . Suppose that for some $\tau, \gamma \geq 0$,*

$$\mathcal{B}(g, \gamma) \cap F_1 \neq \emptyset \text{ and } \inf_{f \in \mathcal{B}(g, \gamma) \cap F_1} \varphi_1(f) \leq \varphi_2(g) + \tau \quad \forall g \in F_2.$$

Then, for any $\delta \in \mathbb{R}$,

$$\text{exs}(\text{lev}_{\leq \delta} \varphi_2, \text{lev}_{\leq \delta + \tau} \varphi_1) \leq \gamma. \quad (7)$$

If in addition

$$\mathcal{B}(g, \gamma) \cap F_2 \neq \emptyset \text{ and } \inf_{f \in \mathcal{B}(g, \gamma) \cap F_2} \varphi_2(f) \leq \varphi_1(g) + \tau \quad \forall g \in F_1,$$

then for any $\varepsilon \geq 0$,

$$\text{exs}(\varepsilon\text{-argmin}_{f \in F_1} \varphi_1(f), (\varepsilon + 2\tau)\text{-argmin}_{f \in F_2} \varphi_2(f)) \leq \gamma. \quad (8)$$

Proof. Let $g \in \text{lev}_{\leq \delta} \varphi_2$. Since φ_1 is lsc and $\mathcal{B}(g, \gamma)$ is compact, there exists $f^* \in \mathcal{B}(g, \gamma) \cap F_1$ such that

$$\varphi_1(f^*) = \inf_{f \in \mathcal{B}(g, \gamma) \cap F_1} \varphi_1(f) \leq \varphi_2(g) + \tau \leq \delta + \tau.$$

We have established that $f^* \in \text{lev}_{\leq \delta + \tau} \varphi_1$. Thus, $\text{dist}(g, \text{lev}_{\leq \delta + \tau} \varphi_1) \leq \gamma$ and (7) follows.

For (8), we note that there exists $f^* \in \text{argmin}_{f \in F_2} \varphi_2(f)$ because F_2 is totally bounded. Thus,

$$\inf_{f \in F_1} \varphi_1(f) \leq \inf_{f \in \mathcal{B}(f^*, \gamma) \cap F_1} \varphi_1(f) \leq \varphi_2(f^*) + \tau = \inf_{f \in F_2} \varphi_2(f) + \tau.$$

Suppose that $g \in \varepsilon\text{-argmin}_{f \in F_1} \varphi_1(f)$. Again, there exists $f^{**} \in \mathcal{B}(g, \gamma) \cap F_2$ such that

$$\varphi_2(f^{**}) = \inf_{f \in \mathcal{B}(g, \gamma) \cap F_2} \varphi_2(f) \leq \varphi_1(g) + \tau \leq \inf_{f \in F_1} \varphi_1(f) + \varepsilon + \tau \leq \inf_{f \in F_2} \varphi_2(f) + \varepsilon + 2\tau.$$

We have established that $f^{**} \in (\varepsilon + 2\tau)\text{-argmin}_{f \in F_2} \varphi_2(f)$. Thus, $\text{dist}(g, (\varepsilon + 2\tau)\text{-argmin}_{f \in F_2} \varphi_2(f)) \leq \gamma$ and (8) follows. \square

We observe that if φ_1 and φ_2 in the proposition are pointwise within δ of each other uniformly on F , then τ can be set to δ and γ to zero. However, the focus on uniform bounds is limiting as it rules out discontinuous functions and especially cases with $F_1 \neq F_2$.

5.1 Confidence Regions

We are then in a position to state the first of the two main results in this section.

5.3 Theorem (confidence region). *For a complete probability space $(\Xi, \mathcal{A}, \mathbb{P})$ and a closed and bounded set $F \subset \text{usc-fcns}(S)$, suppose that $\psi : \Xi \times F \rightarrow (-\infty, \infty]$ is an inf-integrable random lsc function, ξ^1, ξ^2, \dots are independent random elements, each with values in Ξ and distributed according to \mathbb{P} , and $\psi(\xi^1, f)$ is sub-exponential⁷ for all $f \in F$. Given $\alpha \in (0, 1)$ and $\delta > \inf_{f \in F} \mathbb{E}[\psi(\xi^1, f)]$, there exist $\bar{\nu} \in \mathbb{N}$ and $c \in [0, \infty)$ such that for all $\nu \geq \bar{\nu}$*

$$P^\nu \left[\text{exs} \left(\text{argmin}_{f \in F} \mathbb{E}[\psi(\xi^1, f)], \text{lev}_{\leq \delta} \left\{ \frac{1}{\nu} \sum_{j=1}^{\nu} \psi(\xi^j, \cdot) \right\} \right) \leq \frac{c(\log \nu)^{1+1/n}}{\nu^{1/n}} \right] \geq 1 - \alpha.$$

Proof. Let $\varphi : F \rightarrow \mathbb{R}$ have values $\varphi(f) = \mathbb{E}[\psi(\xi^1, f)]$, which is well-defined, lsc, and indeed finite valued due the sub-exponential assumption. Let $\varphi^\nu : \Xi^\nu \times F \rightarrow (-\infty, \infty]$ have values $\varphi^\nu((\xi^1, \dots, \xi^\nu), f) = \nu^{-1} \sum_{j=1}^{\nu} \psi(\xi^j, f)$, which then is a random lsc function on the product probability space. At a given $f \in F$, with probability one, $\varphi^\nu((\xi^1, \dots, \xi^\nu), f) < \infty$ because otherwise $\mathbb{E}[\psi(\xi^1, f)]$ would not have been finite.

As in Theorem 4.3, there is a finite number $N = N(F, \gamma_1)$ of closed balls in $(\text{usc-fcns}(S), d)$ with radius $\gamma_1 > 0$ and center $f_k \in \text{usc-fcns}(S)$ covering F . Without loss of generality, we can assume that $\mathcal{B}(f_k, \gamma_1) \cap F \neq \emptyset$. Moreover, let $f_k^* \in \text{argmin}_{f \in \mathcal{B}(f_k, \gamma_1) \cap F} \varphi(f)$. Since $\psi(\xi^1, f_k^*)$ is sub-exponential, there exists by Bernstein's inequality $\gamma_2 \in (0, \delta - \inf_{f \in F} \varphi(f))$ and $c_0 > 0$ such that

$$P^\nu (|\varphi^\nu((\xi^1, \dots, \xi^\nu), f_k^*) - \varphi(f_k^*)| \geq \gamma_2) \leq 2e^{-\nu c_0 \gamma_2^2} \text{ for all } k = 1, \dots, N.$$

Consequently, as long as

$$2Ne^{-\nu c_0 \gamma_2^2} \leq \alpha \text{ or, equivalently, } \nu \geq \frac{\log N - \log(\alpha/2)}{c_0 \gamma_2^2} \quad (9)$$

we have that

$$P^\nu \left(\max_{k=1, \dots, N} \left| \varphi^\nu((\xi^1, \dots, \xi^\nu), f_k^*) - \varphi(f_k^*) \right| \geq \gamma_2 \right) \leq \alpha.$$

Suppose that we have an event $(\xi^1, \dots, \xi^\nu) \in \Xi^\nu$ where

$$\max_{k=1, \dots, N} \left| \varphi^\nu((\xi^1, \dots, \xi^\nu), f_k^*) - \varphi(f_k^*) \right| < \gamma_2.$$

Next, we apply Proposition 5.2 and start by establishing the required condition. Let $g \in F$ and $\delta_0 = \inf_{f \in F} \varphi(f)$, which is finite because F is compact. Then, there exists $k^* \in \{1, \dots, N\}$ such that $g \in \mathcal{B}(f_{k^*}, \gamma_1)$ and

$$\inf_{f \in \mathcal{B}(g, 2\gamma_1) \cap F} \varphi^\nu((\xi^1, \dots, \xi^\nu), f) \leq \varphi^\nu((\xi^1, \dots, \xi^\nu), f_{k^*}^*) \leq \varphi(f_{k^*}^*) + \gamma_2 \leq \varphi(g) + \delta - \delta_0.$$

⁷A random variable Y is sub-exponential if for some $\lambda \geq 0$, $\mathbb{E}[\exp(\tau(Y - \mathbb{E}Y))] \leq \exp(\tau^2 \lambda^2 / 2)$ for all $|\tau| \leq 1/\lambda$. Another assumption that ensures a Bernstein-type large-deviation result could have been substituted here.

Thus, the first condition in Proposition 5.2 holds with $\gamma = 2\gamma_1$ and $\tau = \delta - \delta_0$, and

$$\text{exs} \left(\text{lev}_{\leq \delta_0} \varphi, \text{lev}_{\leq \delta_0 + \tau} \varphi^\nu((\xi^1, \dots, \xi^\nu), \cdot) \right) \leq 2\gamma_1.$$

Equivalently,

$$\text{exs} \left(\text{argmin}_{f \in F} \varphi(f), \text{lev}_{\leq \delta} \varphi^\nu((\xi^1, \dots, \xi^\nu), \cdot) \right) \leq 2\gamma_1.$$

By Theorem 4.3, there exists $\bar{\varepsilon} > 0$ such that $\log N$ is bounded from above by a term proportional to $\gamma_1^{-n}(\log \gamma_1^{-1})^{n+1}$ for all $\gamma_1 \in (0, \bar{\varepsilon}]$. Thus, there is a constant $c_1 > 0$ such that

$$\frac{\log N - \log(\alpha/2)}{c_0 \gamma_2^2} \leq c_1 \gamma_1^{-n} (\log \gamma_1^{-1})^{n+1} \text{ for } \gamma_1 \in (0, \bar{\varepsilon}]. \quad (10)$$

In view of (9), the right-hand side of (10) provides the rate of increase in sample size that is needed to guarantee an excess of at most $2\gamma_1$ with confidence level $1 - \alpha$. Inverting the expression, we find that γ_1 can be propositional to $\nu^{-1/n}(\log \nu)^{1+1/n}$ as long as ν is sufficiently large, which establishes the conclusion. \square

When $f^* \in \text{argmin}_{f \in F} \mathbb{E}[\psi(\xi^1, f)]$, the theorem guarantees that with probability $1 - \alpha$

$$\text{dist} \left(f^*, \text{lev}_{\leq \delta} \left\{ \frac{1}{\nu} \sum_{j=1}^{\nu} \psi(\xi^j, \cdot) \right\} \right) \leq \frac{c(\log \nu)^{1+1/n}}{\nu^{1/n}}$$

for sufficiently large ν . Hence, the minimizer f^* of (FIP-U) is covered by the given level set when appropriately enlarged with a quantity that vanishes with increasing sample size at nearly the rate $\nu^{-1/n}$. The confidence region is not given in terms of minimizers of the approximating problem (FIP-U) $^\nu$, but rather certain level sets. Membership in such a level set is trivially assessed, does not require solving the approximating problem, and can be used to rule out the optimality of a candidate f . In general, minimizers of (FIP-U) $^\nu$ are not well behaved and depend on the conditioning of (FIP-U) as discussed in [33]. Theorem 5.3 bypasses this issue by considering level sets. Other strengths of Theorem 5.3 are its mild assumption on the (random) objective function ψ and the wide range of constraints that is permitted; the family F can be any bounded closed set in $(\text{usc-fcns}(S), d)$. The functions $f \mapsto \psi(\xi, f)$ is only required to be lsc. The assumption about sub-exponential distribution of $\psi(\xi^1, f)$ can be checked pointwise for each $f \in F$. Actually, this assumption can be relaxed because the proof of Theorem 5.3 only requires that sample averages are sufficiently low relative to the actual expectations, but this merely improves c in the theorem and we omit this refinement.

The practical construction of confidence regions is hampered by the unknown and hard-to-estimate constants c and $\bar{\nu}$ in Theorem 5.3. In practice, coverage may therefore only be guaranteed asymptotically. The other unknown parameter δ is easy to estimate conservatively because for any $f \in F$, the sample average $\frac{1}{\nu} \sum_{j=1}^{\nu} \psi(\xi^j, f)$, using a different sample, furnishes an estimator of $\mathbb{E}[\psi(\xi^1, f)]$, which in turn is an upper bound on $\inf_{f \in F} \mathbb{E}[\psi(\xi^1, f)]$. An effort to select a low δ would obviously result in a smaller level set, but typically also large c and $\bar{\nu}$.

The effect of n on the rate of convergence is profound and in line with the growth of the covering numbers as n increase. It highlights, for example, the fundamental challenge associated with high-dimensional nonparametric estimation already well documented (see [1, 22]). On the positive note, if $n = 1$, which already captures many interesting applications [37], then the convergence rate is nearly ν^{-1} and therefore *faster* than the canonical $\nu^{-1/2}$ rate. If F is restricted to some finite-dimensional subset of $\text{usc-fcns}(S)$, then the covering numbers from Theorem 4.3 can be replaced by much improved ones, typically of order $O(\varepsilon^{-1})$ so that their logarithm is of order $O(\log \varepsilon^{-1})$ and the rate improves from essentially $\nu^{-1/n}$ to $e^{-\nu}$ in Theorem 5.3.

We illustrate the application of Theorem 5.3 on stochastic optimization problems arising in non-parametric statistics.

Example 1: Maximum Likelihood Estimation of Probability Densities. Suppose that we would like to estimate an unknown probability density function $f^0 \in \text{usc-fcns}(S)$. Since we permit densities to have value zero on a subset of S , there is no requirement that the support of f^0 is known; S just needs to contain the support. Given a sample ξ^1, \dots, ξ^ν , which in this case takes values in S , i.e., $\Xi = S$, a maximum likelihood estimator of f^0 over a class $F \subset \text{usc-fcns}(S)$ is any minimizer of

$$\min_{f \in F} -\frac{1}{\nu} \sum_{j=1}^{\nu} \log f(\xi^j)$$

and, in the notation above, $\psi(\xi, f) = -\log f(\xi)$. The function $(\xi, f) \mapsto -\log f(\xi)$ is a random lsc function on the probability space (S, \mathcal{B}, P) , where P is the probability distribution of f^0 and \mathcal{B} contains the Borel sets of $(S, \|\cdot - \cdot\|_\infty)$ supplemented with the necessary probability-zero sets to make the probability space complete. This fact is easily realized because the function is actually lsc jointly in its arguments; see [38] for details.

In this case, $\psi(\xi^1, f)$ being sub-exponential amounts to having F consist of sub-exponential densities. The requirement about inf-integrability is extensively discussed in [38]. For example, suppose F is a nonempty closed subset of

$$\left\{ f \in \text{usc-fcns}(S) \mid \int f(x)dx = 1, \int xf(x)dx \in C, u(x) \leq f(x) \leq v(x), \forall x \in S \right\},$$

where $C \subset \mathbb{R}^n$ is closed and $u, v : S \rightarrow (0, \infty)$, with $v \in \text{usc-fcns}(S)$. Moreover, suppose that the actual density $f^0 \in F$ and for some $\gamma_1 \geq 0, \gamma_2 > 0$, and, $\zeta_1, \zeta_2 \in \mathbb{R}$,

$$u(x) \geq e^{-\gamma_1 \|x\|_\infty + \zeta_1} \quad \text{and} \quad v(x) \leq e^{-\gamma_2 \|x\|_\infty + \zeta_2}.$$

All the assumptions of Theorem 5.3 are then satisfied. The requirement that F is bounded is automatically satisfied because $f \geq 0$ for all $f \in F$.

In this example the maximum likelihood estimator finds the best estimate that satisfies the given pointwise bounds and moment restriction. There is *no* requirement that the actual density or its estimate should be smooth or even continuous. Of course, a large variety of other constraints can be

brought in too; see [38] for some possibilities.

We recall that a subset $F_0 \subset \text{usc-fcns}(S)$ is *equi-usc* [31, Sect. 7.B] if there exists $\delta : S \times (0, \infty) \times (0, \infty) \rightarrow (0, \infty)$ such that for any $\varepsilon, \rho > 0$, $\bar{x} \in S$, and $f \in F_0$,

$$\sup_{x \in B_\infty(\bar{x}, \delta(\bar{x}, \varepsilon, \rho))} f(\bar{x}) \leq \max\{f(x) + \varepsilon, -\rho\}.$$

If F_0 is a singleton, then the condition reduces to that of usc. If F_0 contains only Lipschitz continuous functions, or only piecewise Lipschitz continuous functions, or only finite-valued concave functions on \mathbb{R}^n , to mention some examples, then F_0 is equi-usc.

Example 2: Least-Squares Regression. Suppose that we are given the random design model

$$\mathbf{y}^j = f^0(\mathbf{x}^j) + \mathbf{z}^j, \quad j = 1, 2, \dots, \nu$$

where $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^\nu$ are independent and identically distributed n -dimensional random vectors that take values in a closed set $S \subset \mathbb{R}^n$, $\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^\nu$ are zero-mean random variables that are also independent of $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^\nu$, and $f^0 : S \rightarrow \mathbb{R}$ is an unknown function to be estimated based on observations of $\boldsymbol{\xi}^1 = (\mathbf{x}^1, \mathbf{y}^1)$. In this case, $\Xi = S \times \mathbb{R}$, again we adopt a sigma-algebra that contains the Borel sets on Ξ and that results in a complete probability space under the distribution of $(\mathbf{x}^1, \mathbf{y}^1)$. The least-squares estimator of f^0 over the class $F \subset \text{usc-fcns}(S)$ is then any minimizer of

$$\min_{f \in F} \frac{1}{\nu} \sum_{j=1}^{\nu} (\mathbf{y}^j - f(\mathbf{x}^j))^2.$$

Resulting estimates furnish approximations of f^0 that in an engineering design context can be maximized to find an optimal design without any (additional) costly simulation of system performance. The only simulations required are those needed to generate a data set $\{(x^j, y^j), j = 1, \dots, \nu\}$.

In this case, $\psi((x, y), f) = (y - f(x))^2$. Since $(x, f) \mapsto f(x)$ is usc and thus measurable, we also have that ψ is measurable. Consequently, ψ is a random lsc function provided that F is equi-usc, an assumption that provides the necessary pointwise convergence (cf. [31, Thm. 7.10]). Its nonnegativity ensures that ψ is also locally inf-integrable.

A confidence region for $f^0 \in F$ emerges from Theorem 5.3 when $(\mathbf{y}^1 - f(\mathbf{x}^1))^2$ is sub-exponential for all $f \in F$. For example, this will be the case when \mathbf{z}^1 and every component of \mathbf{x}^1 are sub-Gaussian, and for some $\gamma, \zeta \in \mathbb{R}$,

$$f \in F \implies |f(x)| \leq \gamma \|x\|_\infty + \zeta, \quad \forall x \in S.$$

Since f^0 must be a minimizer of $\min_{f \in F} \mathbb{E}[(\mathbf{y}^1 - f(\mathbf{x}^1))^2]$, provided that $f^0 \in F$, Theorem 5.3 guarantees that f^0 is covered by the stipulated level set when appropriately enlarged.

5.2 Rates of Convergence under Hölder Condition

Theorem 5.3 does not rule out the possibility that the limit of the given level sets strictly contains $\text{argmin}_{f \in F} \mathbb{E}[\psi(\boldsymbol{\xi}^1, f)]$. In fact, this cannot be ruled out unless additional assumptions are brought in;

[33] contains a discussion. Still, a Hölder condition enables us to “reverse” Theorem 5.3 and quantify the rate of convergence of the excess of minimizers of $(\text{FIP-U})^\nu$ over those of (FIP-U) . Since it is relatively straightforward, we also address approximating constraints. Although the approximating constraints can be rather general, the rate of convergence in the following theorem depends on the rate with which the approximating feasible set approaches the actual one. Thus, it is not immediately clear how the piecewise affine functions discussed in Section 3, which have unknown rate of convergence, can be used for constructing these approximations.

As usual, we let $\alpha^\nu = o(r^\nu)$ imply that for every $\delta > 0$ there exists $\bar{\nu}$ such that $\alpha^\nu \leq \delta r^\nu$ for all $\nu \geq \bar{\nu}$.

5.4 Theorem (rate of convergence). *For a complete probability space $(\Xi, \mathcal{A}, \mathbb{P})$ and closed and bounded sets $F^\nu, F^0 \subset F \subset \text{usc-fcns}(S)$, suppose that $\psi : \Xi \times F \rightarrow (-\infty, \infty]$ is a random lsc function for which there exist $p \in (0, \infty)$ and integrable random variable $\kappa : \Xi \rightarrow [0, \infty)$ such that*

$$|\psi(\xi, f) - \psi(\xi, g)| \leq \kappa(\xi)[\mathcal{d}(f, g)]^p \text{ for all } f, g \in F \text{ and } \xi \in \Xi.$$

Suppose also that ξ^1, ξ^2, \dots are independent random elements, each with values in Ξ and distributed according to \mathbb{P} , and $\psi(\xi^1, f)$ is sub-exponential for all $f \in F$. Let

$$r^\nu = \nu^{\frac{-1}{2+n/p}} (\log \nu)^{\frac{1+n}{2+n/p}}.$$

If $\mathbb{H}(F^\nu, F^0) = o(\min\{r^\nu, (r^\nu)^{1/p}\})$ and $\alpha \in (0, 1)$, then there exist $c \in [0, \infty)$ and $\bar{\nu} \in \mathbb{N}$ such that for $\nu \geq \bar{\nu}$ and $\varepsilon^\nu \geq 0$,

$$P^\nu \left[\text{exs} \left(\varepsilon^\nu - \text{argmin}_{f \in F^\nu} \frac{1}{\nu} \sum_{j=1}^{\nu} \psi(\xi^j, f), (\varepsilon^\nu + cr^\nu) - \text{argmin}_{f \in F^0} \mathbb{E}[\psi(\xi^1, f)] \right) \leq \mathbb{H}(F^\nu, F^0) \right] \geq 1 - \alpha.$$

Proof. Let $\zeta > 0$. Since $\kappa(\xi^1)$ is integrable, there exists $\bar{\nu}_0 \in \mathbb{N}$ such that $P^\nu(|\nu^{-1} \sum_{j=1}^{\nu} \kappa(\xi^j) - \mathbb{E}[\kappa(\xi^1)]| \geq \zeta) \leq \alpha/2$ for all $\nu \geq \bar{\nu}_0$. Let $\gamma_1 > 0$. As in Theorem 4.3, there is a finite number $N = N(F, \gamma_1/2)$ of closed balls in $(\text{usc-fcns}(S), \mathcal{d})$ with radius $\gamma_1/2$ and center f'_k covering F . To make sure that the balls are centered at points in F , we can always select some other centers $f_k \in F$ and balls with radius γ_1 and still cover F .

Let φ and φ^ν be as defined in the proof of Theorem 5.3. We note that ψ is locally inf-integrable due to the Hölder condition and the pointwise sub-exponential property. Since $\psi(\xi^1, f_k)$ is sub-exponential, there exists by Bernstein’s inequality $\bar{\gamma}_2 > 0$ and $c_0 > 0$ such that for $\gamma_2 \in [0, \bar{\gamma}_2]$,

$$P^\nu(|\varphi^\nu((\xi^1, \dots, \xi^\nu), f_k) - \varphi(f_k)| \geq \gamma_2) \leq 2e^{-\nu c_0 \gamma_2^2} \text{ for all } k = 1, \dots, N.$$

Consequently, as long as $\nu \geq \bar{\nu}_0$ and $2Ne^{-\nu c_0 \gamma_2^2} \leq \alpha/2$, or, equivalently,

$$\nu \geq \max \left\{ \bar{\nu}_0, \frac{\log N - \log(\alpha/4)}{c_0 \gamma_2^2} \right\}$$

we have that

$$P^\nu \left(\max_{k=1, \dots, N} \left| \varphi^\nu((\xi^1, \dots, \xi^\nu), f_k) - \varphi(f_k) \right| \geq \gamma_2 \text{ or } \left| \frac{1}{\nu} \sum_{j=1}^\nu \kappa(\xi^j) - \mathbb{E}[\kappa(\xi^1)] \right| \geq \zeta \right) \leq \alpha.$$

Suppose that we have an event $(\xi^1, \dots, \xi^\nu) \in \Xi^\nu$ where

$$\max_{k=1, \dots, N} \left| \varphi^\nu((\xi^1, \dots, \xi^\nu), f_k) - \varphi(f_k) \right| < \gamma_2 \text{ and } \left| \frac{1}{\nu} \sum_{j=1}^\nu \kappa(\xi^j) - \mathbb{E}[\kappa(\xi^1)] \right| < \zeta.$$

Next, we apply Proposition 5.2 for the lsc functions $\bar{\varphi} : F^0 \rightarrow \mathbb{R}$ given by $\bar{\varphi}(f) = \varphi(f)$ and $\bar{\varphi}^\nu : F^\nu \rightarrow \mathbb{R}$ given by $\bar{\varphi}^\nu(f) = \varphi^\nu((\xi^1, \dots, \xi^\nu), f)$. In view of the Hölder assumption on ψ , this implies that $\bar{\varphi}^\nu$ is finite when defined. Moreover, for all $f, g \in F$,

$$|\varphi(f) - \varphi(g)| \leq \mathbb{E}[\kappa(\xi^1)] [d(f, g)]^p.$$

Let $\delta^\nu = \mathbb{H}(F^\nu, F^0)$. Suppose that $f \in F^\nu$. Then, there is $f' \in F^0$ and $k^* \in \{1, \dots, N\}$ such that $d(f, f') \leq \delta^\nu$ and $d(f', f_{k^*}) \leq \gamma_1$. Thus,

$$\begin{aligned} \inf_{g \in B(f, \delta^\nu) \cap F} \bar{\varphi}(g) &\leq \bar{\varphi}(f') \leq \bar{\varphi}(f_{k^*}) + \mathbb{E}[\kappa(\xi^1)] \gamma_1^p \\ &< \bar{\varphi}^\nu(f_{k^*}) + \gamma_2 + \mathbb{E}[\kappa(\xi^1)] \gamma_1^p \\ &\leq \bar{\varphi}^\nu(f) + \gamma_2 + \mathbb{E}[\kappa(\xi^1)] \gamma_1^p + \frac{1}{\nu} \sum_{j=1}^\nu \kappa(\xi^j) (\delta^\nu + \gamma_1)^p \\ &\leq \bar{\varphi}^\nu(f) + \gamma_2 + \mathbb{E}[\kappa(\xi^1)] (\gamma_1^p + (\delta^\nu + \gamma_1)^p) + \zeta (\delta^\nu + \gamma_1)^p. \end{aligned}$$

Similarly, suppose that $f \in F^0$. Then, there is $f' \in F^\nu$ and $k^* \in \{1, \dots, N\}$ such that $d(f, f') \leq \delta^\nu$ and $d(f', f_{k^*}) \leq \gamma_1$. Consequently,

$$\begin{aligned} \inf_{g \in B(f, \delta^\nu) \cap F} \bar{\varphi}^\nu(g) &\leq \bar{\varphi}^\nu(f') \leq \bar{\varphi}^\nu(f_{k^*}) + \frac{1}{\nu} \sum_{j=1}^\nu \kappa(\xi^j) \gamma_1^p \\ &< \bar{\varphi}(f_{k^*}) + \gamma_2 + \frac{1}{\nu} \sum_{j=1}^\nu \kappa(\xi^j) \gamma_1^p \\ &\leq \bar{\varphi}(f) + \mathbb{E}[\kappa(\xi^1)] (\delta^\nu + \gamma_1)^p + \gamma_2 + \frac{1}{\nu} \sum_{j=1}^\nu \kappa(\xi^j) \gamma_1^p \\ &\leq \bar{\varphi}(f) + \mathbb{E}[\kappa(\xi^1)] [(\delta^\nu + \gamma_1)^p + \gamma_1^p] + \gamma_2 + \zeta \gamma_1^p. \end{aligned}$$

Thus, we have shown that the conditions of Proposition 5.2 hold for the functions $\bar{\varphi}$ and $\bar{\varphi}^\nu$ with

$$\delta^\nu \text{ and } \tau_0 = \gamma_2 + \mathbb{E}[\kappa(\xi^1)] (\gamma_1^p + (\delta^\nu + \gamma_1)^p) + \zeta (\delta^\nu + \gamma_1)^p$$

as the two error parameters (γ and τ) and we therefore have that

$$\text{exs}(\varepsilon^\nu\text{-argmin}_{f \in F^\nu} \varphi^\nu((\xi^1, \dots, \xi^\nu), f), (\varepsilon^\nu + 2\tau_0)\text{-argmin}_{f \in F^0} \varphi(f)) \leq \delta^\nu.$$

By Theorem 4.3, $\log N$ is bounded from by a term proportional to $\gamma_1^{-n}(\log \gamma_1^{-1})^{n+1}$ for sufficiently small γ_1 . Thus, there exist constants $c_1, c_2 > 0$ such that

$$\frac{\log N - \log(\alpha/4)}{c_0\gamma_2^2} \leq c_1\gamma_1^{-n}(\log \gamma_1^{-1})^{n+1}\gamma_2^{-2} + c_2\gamma_2^{-2},$$

which gives the rate of growth in ν as γ_1 and γ_2 vanish. For $\tau > 0$, the error τ_0 can be kept below τ if γ_1 is proportional to $\tau^{1/p}$, γ_2 is proportional to τ , δ^ν is proportional to $\min\{\tau, \tau^{1/p}\}$, and the (positive) proportionality constants are selected sufficiently close to zero. In view of these choices about γ_1 and γ_2 , there is a constant $c_3 > 0$ such that

$$c_1\gamma_1^{-n}(\log \gamma_1^{-1})^{n+1}\gamma_2^{-2} + c_2\gamma_2^{-2} \leq c_3\tau^{-2-n/p}(\log \tau^{-1/p})^{n+1}.$$

With ν above $\bar{\nu}_0$ as well as the previous right-hand side, or equivalently for some $c_4 > 0$,

$$\tau \geq c_4\nu^{\frac{-1}{2+n/p}}(\log \nu)^{\frac{1+n}{2+n/p}},$$

we ensure the required confidence level and the conclusion follows. \square

A corollary of the theorem for the case with $\varepsilon^\nu = 0$ and $F^\nu = F^0 = F$ is that

$$\operatorname{argmin}_{f \in F} \frac{1}{\nu} \sum_{j=1}^{\nu} \psi(\xi^j, f) \subset c r^\nu \operatorname{argmin}_{f \in F} \mathbb{E}[\psi(\xi^1, f)]$$

with at least probability $1 - \alpha$. Thus, minimizers of $(\text{FIP-U})^\nu$ converge at the rate r^ν to a minimizer of (FIP-U) . The rate depends on the Hölder coefficient p as well as the dimension n of the space of function under considerations.

We illustrate the assumptions of the theorem for two stochastic optimization problems arising in nonparametric statistics, but start with an intermediate result.

5.5 Proposition *For Lipschitz continuous functions $f, g \in \text{usc-fcns}(S)$ with common modulus $\kappa \in [0, \infty)$,*

$$|f(x) - g(x)| \leq (1 + \kappa)e^{\rho(x)} \mathbf{d}(f, g) \text{ for all } x \in S,$$

where $\rho(x) = \max\{\|x\|_\infty, |f(x)|, |g(x)|\}$.

Proof. Let $x \in S$. The first result is trivial if $\rho(x) = \infty$. Suppose that $\rho(x) < \infty$. From Lemma 4.1, $\mathbf{d}(f, g) \geq e^{-\rho(x)} \hat{\mathbf{d}}_{\rho(x)}(f, g)$. Set $\tau \in (\hat{\mathbf{d}}_{\rho(x)}(f, g), \infty)$. Again, by Lemma 4.1, there exists $y \in \mathcal{B}_\infty(x, \tau)$ such that $f(y) \geq g(x) - \tau$. Thus, $g(x) - f(x) = g(x) - f(y) + f(y) - f(x) \leq \tau + \kappa\tau$. A similar argument establishes that $f(x) - g(x) \leq \tau + \kappa\tau$. Hence, by letting τ tends to its lower limit, we obtain that $|f(x) - g(x)| \leq (1 + \kappa)\hat{\mathbf{d}}_{\rho(x)}(f, g)$ and the conclusion follows. \square

Example 3: Least-Squares Regression. We return to the setting of Example 2, but now let F be a family that contains only Lipschitz continuous functions with common modulus $\kappa_0 \geq 0$. Suppose

also that \mathbf{z}^1 and every component of \mathbf{x}^1 are sub-Gaussian, the unknown function $f^0 \in F$, and there exists $\beta < \infty$ such that $f(0) \leq \beta$ for all $f \in F$. Then, F is equi-usc and there are $\gamma, \zeta \in \mathbb{R}$ such that $|f(x)| \leq \gamma \|x\|_\infty + \zeta$ for all $f \in F$. Proposition 5.5 then ensures that the Hölder condition in Theorem 5.4 holds with $p = 2$ and $\kappa((x, y)) = (1 + \kappa_0)^2 \exp(2 \max\{\|x\|_\infty, \gamma \|x\|_\infty + \zeta\})$, which is integrable in view of the sub-Gaussianity of \mathbf{x}^1 .

Using the bound on $|f(x)|$, we also have that $(\mathbf{y}^1 - f(\mathbf{x}^1))^2 = (f^0(\mathbf{x}^1) - f(\mathbf{x}^1) + \mathbf{z}^1)^2$ is sub-exponential. The assumptions of Theorem 5.4 therefore hold,

$$r^\nu = \nu^{\frac{-2}{4+n}} (\log \nu)^{\frac{1+n}{2+n/2}},$$

and, for closed $F^\nu, F^0 \subset F$, there exist $c \in [0, \infty)$ and $\bar{\nu} \in \mathbb{N}$ such that

$$P^\nu \left[\text{exs} \left(\operatorname{argmin}_{f \in F^\nu} \frac{1}{\nu} \sum_{j=1}^{\nu} (\mathbf{y}^j - f(\mathbf{x}^j))^2, \quad cr^\nu - \operatorname{argmin}_{f \in F^0} \mathbb{E}[(\mathbf{y}^1 - f(\mathbf{x}^1))^2] \right) \leq \mathbb{H}(F^\nu, F^0) \right] \geq 1 - \alpha$$

provided that $\nu \geq \bar{\nu}$ and $\mathbb{H}(F^\nu, F^0) = o(r^\nu)$. Thus, when $\mathbb{H}(F^\nu, F^0) = 0$ and $\hat{f}^\nu \in \operatorname{argmin}_{f \in F^\nu} \frac{1}{\nu} \sum_{j=1}^{\nu} (\mathbf{y}^j - f(\mathbf{x}^j))^2$ is measurable,

$$P^\nu \left(\mathbb{E}[(\hat{f}^\nu(\mathbf{x}^1) - f^0(\mathbf{x}^1))^2] \leq cr^\nu \right) \geq 1 - \alpha.$$

The rates developed here apply in rather general settings and remain in effect even if $f^0 \notin F^0$. More specific settings give improved results as in the case of regression with fixed design and Lipschitz continuous functions defined on compact convex subset [44, p. 333] and in the univariate case [15].

Example 4: Least-Squares Probability Density Estimation. We return to the setting of Example 1, but now consider the least-squares estimator of f^0 , which is any minimizer of

$$\min_{f \in F^\nu} -\frac{2}{\nu} \sum_{j=1}^{\nu} f(\boldsymbol{\xi}^j) + \int [f(x)]^2 dx.$$

This estimator is motivated by the fact that the unknown function

$$f^0 \in \operatorname{argmin}_{f \in F} \int [f(x) - f^0(x)]^2 dx = \operatorname{argmin}_{f \in F} -2\mathbb{E}[f(\boldsymbol{\xi}^1)] + \int [f(x)]^2 dx$$

whenever $f^0 \in F$. To make the case rather concrete, let $\kappa \in [0, \infty)$ and for some bounded function $h : S \rightarrow [0, \infty)$, with $\int h(x) dx < \infty$,

$$F = \left\{ f \in \text{usc-fcns}(S) \mid \int f(x) dx = 1, \quad 0 \leq f(x) \leq h(x), \quad |f(x) - f(y)| \leq \kappa \|x - y\|_\infty, \quad \forall x, y \in S \right\},$$

which can be shown to be closed and bounded; see arguments in [38]. In this case, $(\boldsymbol{\xi}, x) \mapsto \psi(\boldsymbol{\xi}, f) = -2f(\boldsymbol{\xi}) + \int [f(x)]^2 dx$ is a random lsc function as can be seen by invoking Fatou's Lemma and pointwise convergence; again see [38]. Then, $\psi(\boldsymbol{\xi}^1, f)$ is sub-exponential for all $f \in F$ as it is in fact bounded.

It remains to check the Hölder condition in Theorem 5.4. Suppose that $\mathbb{E}[\exp(\|\xi^1\|_\infty)] < \infty$. In view of Proposition 5.5, if the integral term in ψ had not been present, then the condition holds with $p = 1$; Lipschitz continuity and the fact that $\mathbb{E}[\exp(\|\xi^1\|_\infty)] < \infty$ ensures integrability of the Hölder modulus. If S were compact, then ψ would still satisfy the condition with $p = 1$. For a noncompact S , the argument needs to be slightly modified by first “ignoring” the integral term and second reintroduce it in a slightly generalized version of Theorem 5.4. We omit the details.

In summary, for the given F and under the assumption that $\mathbb{E}[\exp(\|\xi^1\|_\infty)] < \infty$, we can show by invoking Theorem 5.4 (or the mentioned extensions) that for any $\alpha \in (0, 1)$ there exist $c \in [0, \infty)$ and $\bar{\nu} \in \mathbb{N}$ such that for every $\nu \geq \bar{\nu}$ and $\varepsilon^\nu \geq 0$

$$P^\nu \left[\varepsilon^\nu - \operatorname{argmin}_{f \in F} - \frac{2}{\nu} \sum_{j=1}^{\nu} f(\xi^j) + \int [f(x)]^2 dx \subset (\varepsilon^\nu + c r^\nu) - \operatorname{argmin}_{f \in F} - 2\mathbb{E}[f(\xi^1)] + \int [f(x)]^2 dx \right] \geq 1 - \alpha \text{ with } r^\nu = \nu^{\frac{-1}{2+n}} (\log \nu)^{\frac{1+n}{2+n}}.$$

A sharper result is available in the univariate case over the class of nonincreasing convex functions [13].

Acknowledgements. This work is supported in parts by DARPA under grants HR0011-14-1-0060 and HR0011-8-34187, and Office of Naval Research (Science of Autonomy Program) under grant N00014-17-1-2372.

Appendix

Proof of Theorem 4.4. Let $\rho > 0$ and $F = \{f \in \text{usc-fcns}(\mathbb{R}^d) \mid f(x) \geq -\rho \text{ for at least one } x \in [0, \rho]^n\}$. We show that F cannot be covered with a lower number of balls than stipulated. Clearly, $\text{dist}_\infty(0, \text{hypo } f) \leq \rho$ for all $f \in F$. Thus, in view of (3), $\mathfrak{d}(0, f) \leq \rho + 1$ for all $f \in F$, where 0 is the zero function on \mathbb{R}^n , and F is therefore bounded.

Next, let $\varepsilon \in (0, \rho e^{-\rho}/6]$. We discretize $[0, \rho]^n$ by defining $x_i^k = k\rho/\nu_\varepsilon$, $k = 1, \dots, \nu_\varepsilon - 1$ and $i = 1, \dots, n$, where

$$\nu_\varepsilon = \left\lfloor \frac{\rho e^{-\rho}}{3\varepsilon} \right\rfloor \geq 2,$$

with $\lfloor a \rfloor$ being the largest integer not exceeding a . The discretization of $[0, \rho]^n$ then contains the points $(x_1^{k_1}, x_2^{k_2}, \dots, x_n^{k_n})$, with $k_i \in \{1, 2, \dots, \nu_\varepsilon - 1\}$ and $i = 1, \dots, n$. Clearly, the distance between any two such points in the sup-norm is at least $\rho/\nu_\varepsilon \geq 3\varepsilon e^\rho$. We carry out a similar discretization of $[-\rho, 0]$ and define $y^l = l\rho/\nu_\varepsilon$, $l = 1, \dots, \nu_\varepsilon$. The functions that are finite on the discretization points of $[0, \rho]^n$, with values at each such point equal to y^l for some l , and have value minus infinity elsewhere are given by F_ε , i.e.,

$$F_\varepsilon = \{f \in \text{usc-fcns}(\mathbb{R}^n) \mid \text{for each } x = (x_1^{k_1}, \dots, x_n^{k_n}), \text{ with } k_i \in \{1, 2, \dots, \nu_\varepsilon - 1\}, f(x) = y^l \text{ for some } l = 1, \dots, \nu_\varepsilon; f(x) = -\infty \text{ otherwise}\}.$$

Certainly, $F_\varepsilon \subset F$. We next define

$$G_\varepsilon(f) = \{g \in \text{usc-fcns}(\mathbb{R}^n) \mid \hat{d}_\rho(f, g) \leq \varepsilon e^\rho\}, \quad \text{for } f \in \text{usc-fcns}(\mathbb{R}^n).$$

We establish that $G_\varepsilon(f) \cap G_\varepsilon(f') = \emptyset$ for $f, f' \in F_\varepsilon, f \neq f'$. Suppose for the sake of a contradiction that there is a g with $g \in G_\varepsilon(f)$ and $g \in G_\varepsilon(f')$ for $f, f' \in F_\varepsilon, f \neq f'$. Then, $\hat{d}_\rho(f, g) \leq \varepsilon e^\rho$ and $\hat{d}_\rho(f', g) \leq \varepsilon e^\rho$. However, since $f \neq f'$, there exists a point $x \in [0, \rho]^n$ with $|f(x) - f'(x)| \geq 3\varepsilon e^\rho$. Without loss of generality, suppose that $f(x) \geq f'(x) + 3\varepsilon e^\rho$. Since $f(z), f'(z) = -\infty$ for all $z \neq x$ with $\|z - x\|_\infty < 3\varepsilon e^\rho$, we have that $\hat{d}_\rho(f, g) \leq \varepsilon e^\rho$ implies that $g(z) \geq f(x) - \varepsilon e^\rho$ for some $z \in \mathcal{B}(x, \varepsilon e^\rho)$. Moreover, $\hat{d}_\rho(f', g) \leq \varepsilon e^\rho$ implies that $g(z) \leq f'(x) + \varepsilon e^\rho \leq f(x) - 3\varepsilon e^\rho + \varepsilon e^\rho = f(x) - 2\varepsilon e^\rho$ for all $z \in \mathcal{B}(x, \varepsilon e^\rho)$. Since this is not possible for g , we have reached a contradiction. Thus, $G_\varepsilon(f) \cap G_\varepsilon(f') = \emptyset$ for $f, f' \in F_\varepsilon, f \neq f'$.

By Lemma 4.1, for any $f \in \text{usc-fcns}(\mathbb{R}^n)$,

$$d(f, g) \geq e^{-\rho} \hat{d}_\rho(f, g) > e^{-\rho} \varepsilon e^\rho = \varepsilon \text{ for all } g \notin G_\varepsilon(f).$$

Hence, for $f \in F_\varepsilon$, an d -ball with radius ε that contains f needs to be centered at some $g \in G_\varepsilon(f)$. Since the sets $G_\varepsilon(f), f \in F_\varepsilon$, are nonoverlapping, a cover of F_ε by d -balls with radius ε must involve a number of balls that is at least as great as the number of functions in F_ε , which is $\nu_\varepsilon^{m_\varepsilon}$, where $m_\varepsilon = (\nu_\varepsilon - 1)^n$. Thus,

$$\log N(F, \varepsilon) \geq \nu_\varepsilon^n \log \nu_\varepsilon \geq \left(\frac{\rho e^{-\rho}}{3\varepsilon} - 2\right)^n \log \left(\frac{\rho e^{-\rho}}{3\varepsilon} - 1\right). \quad (11)$$

Let $c_1 = |\log(\rho e^{-\rho}/4)|$ and $\bar{\varepsilon} = \min\{\rho e^{-\rho}/12, e^{-2c_1}\}$. Continuing from (11), we then find that

$$\log N(F, \varepsilon) \geq \left(\frac{\rho e^{-\rho}}{6}\right)^n \left[1 + \frac{\log(\rho e^{-\rho}/4)}{\log \varepsilon^{-1}}\right] \frac{1}{\varepsilon^n} \log \frac{1}{\varepsilon}.$$

Since $\log \varepsilon^{-1} \geq 2|\log(\rho e^{-\rho}/4)|$ for $\varepsilon \in (0, \bar{\varepsilon}]$, we have that

$$\log N(F, \varepsilon) \geq \left(\frac{\rho e^{-\rho}}{6}\right)^n \frac{1}{2} \frac{1}{\varepsilon^n} \log \frac{1}{\varepsilon} \quad \text{for } \varepsilon \in (0, \bar{\varepsilon}],$$

and the conclusion is reached. □

References

- [1] F. Balabdaoui and J. A. Wellner. Estimation of a k -monotone density: characterizations, consistency and minimax lower bounds. *Statistica Neerlandica*, 64(1):45–70, 2010.
- [2] D. Bampou and D. Kuhn. Polynomial approximations for continuous linear programs. *SIAM Journal on Optimization*, 22:628–648, 2012.
- [3] P. L. Bartlett, S. R. Kulkarni, and S. E. Posner. Covering numbers for real-valued function classes. *IEEE Transactions on Information Theory*, 43(5):1721–1724, Sep 1997.

- [4] G. Bayraksan and D.P. Morton. Assessing solution quality in stochastic programs. *Mathematical Programming*, 108:495–514, 2006.
- [5] M. S. Birman and M. Z. Solomjak. Piecewise-polynomial approximation of functions of the classes w_p^α . *Mathematics of the USSR-Sbornik*, 73:295–317, 1967.
- [6] E. M. Bronshtein. ε -entropy of convex sets and functions. *Siberian Mathematical Journal*, 17(3):393–398, 1976.
- [7] A. Brudnyi. On covering numbers of sublevel sets of analytic functions. *J. Approximation Theory*, 162(1):72 – 93, 2010.
- [8] Y. Cui, J.-S. Pang, and B. Sen. Composite difference-max programs for modern statistical estimation problems. *ArXiv e-prints*, 2018.
- [9] M. Cule, R.J. Samworth, and M. Stewart. Maximum likelihood estimation of a multi-dimensional log-concave density. *J. Royal Statistical Society Series B*, 72:545–600, 2010.
- [10] O. Devolder, F. Glineur, and Y. Nesterov. Solving infinite-dimensional optimization problems by polynomial approximation. In M. Diehl, F. Glineur, E. Jarlebring, and W. Michiels, editors, *Recent Advances in Optimization and its Applications in Engineering*, pages 31–40. Springer, Berlin, 2010.
- [11] R. M. Dudley. Metric entropy of some classes of sets with differentiable boundaries. *J. Approximation Theory*, 10(3):227–236, 1974.
- [12] A. Georghiou, W. Wiesemann, and D. Kuhn. Generalized decision rule approximations for stochastic programming via liftings. *Mathematical Programming*, 152(1-2):301–338, 2015.
- [13] P. Groeneboom, G. Jongbloed, and J. A. Wellner. Estimation of a convex function: characterizations and asymptotic theory. *Annals of Statistics*, 29:1653–1698, 2001.
- [14] A. Guntuboyina and B. Sen. Covering numbers for convex functions. *IEEE Transactions on Information Theory*, 59(4):1957–1965, 2013.
- [15] A. Guntuboyina and B. Sen. Global risk bounds and adaptation in univariate convex regression. *Probability Theory and Related Fields*, 163:379–411, 2015.
- [16] Y. Guo, P. L. Bartlett, J. Shawe-Taylor, and R. C. Williamson. Covering numbers for support vector machines. *IEEE Transactions on Information Theory*, 48(1):239–250, Jan 2002.
- [17] G. A. Hanasusanto, W. Wiesemann, and D. Kuhn. K-adaptability in two-stage robust binary programming. *Operations Research*, 63(4):877–891, 2015.
- [18] P. Hartman. On functions representable as a difference of convex functions. *Pacific J. Mathematics*, 9:707–713, 1959.

- [19] J.L. Higle and S. Sen. Statistical verification of optimality conditions for stochastic programs with recourse. *Annals of Operations Research*, 30:215–240, 1991.
- [20] J.L. Higle and S. Sen. Duality and statistical tests of optimality for two stage stochastic programs. *Mathematical Programming*, 75:257–275, 1996.
- [21] R. Horst and N. V. Thoai. DC programming: Overview. *J. Optimization Theory and Applications*, 103(1):1–43, 1999.
- [22] A. K. H. Kim and R. J. Samworth. Global rates of convergence in log-concave density estimation. *Annals of Statistics*, 44:2756–2779, 2016.
- [23] A. N. Kolmogorov and V. M. Tikhomirov. Epsilon-entropy and epsilon-capacity of sets in functional spaces. *American Mathematical Society Translations, Series 2*, 17:277–364, 1961.
- [24] T. Kühn. Covering numbers of Gaussian reproducing kernel Hilbert spaces. *Journal of Complexity*, 27(5):489–499, 2011.
- [25] M. Lamm and S. Lu. Generalized conditioning based approaches to computing confidence intervals for solutions to stochastic variational inequalities. *Mathematical Programming B*, to appear, 2018.
- [26] S. Lu, Y. Liu, L. Yin, and K. Zhang. Confidence intervals and regions for the lasso by using stochastic variational inequality techniques in optimization. *J. Royal Statistical Society: Series B*, 79:589–611, 2017.
- [27] W. K. Mak, D. P. Morton, and R. K. Wood. Monte Carlo bounding techniques for determining solution quality in stochastic programs. *Operations Research Letters*, 24:47–56, 1999.
- [28] M. Miller. Binary classification using piecewise affine functions. Master’s thesis, Naval Postgraduate School, Monterey, California, June 2019.
- [29] V.I. Norkin, G.C. Pflug, and A. Ruszczyński. A branch and bound method for stochastic global optimization. *Mathematical Programming*, 83:425–450, 1998.
- [30] M. Pontil. A note on different covering numbers in learning theory. *J. Complexity*, 19(5):665–671, 2003.
- [31] R.T. Rockafellar and R. J-B Wets. *Variational Analysis*, volume 317 of *Grundlehren der Mathematischen Wissenschaft*. Springer, 3rd printing-2009 edition, 1998.
- [32] J. O. Royset. Optimality functions in stochastic programming. *Mathematical Programming*, 135(1):293–321, 2012.
- [33] J. O. Royset. Approximations and solution estimates in optimization. *Mathematical Programming*, 170(2):479–506, 2018.

- [34] J. O. Royset and R. J-B Wets. From data to assessments and decisions: Epi-spline technology. In A. Newman, editor, *INFORMS Tutorials*. INFORMS, Catonsville, 2014.
- [35] J. O. Royset and R. J-B Wets. Multivariate epi-splines and evolving function identification problems. *Set-Valued and Variational Analysis*, 24(4):517–545, 2016. Erratum: pp. 547-549.
- [36] J. O. Royset and R. J-B Wets. Variational theory for optimization under stochastic ambiguity. *SIAM J. Optimization*, 27(2):1118–1149, 2017.
- [37] J. O. Royset and R. J-B Wets. On univariate function identification problems. *Mathematical Programming B*, 168(1-2):449–474, 2018.
- [38] J. O. Royset and R. J-B Wets. Variational analysis of constrained M-estimators. *ArXiv e-prints*, 2018.
- [39] G. Salinetti and R. J-B Wets. On the convergence in distribution of measurable multifunctions (random sets), normal integrands, stochastic processes and stochastic infima. *Mathematics of Operations Research*, 11(3):385–419, 1986.
- [40] G. Salinetti and R. J-B Wets. On the hypo-convergence of probability measures. In *Optimization and Related Fields, Proc., Erice 1984, Lecture Notes in Mathematics 1190*, pages 371–395. Springer, 1986.
- [41] E. Seijo and B. Sen. Nonparametric least squares estimation of a multivariate convex regression. *Annals of Statistics*, 39:1633–1657, 2011.
- [42] A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory*. SIAM, 2. edition, 2014.
- [43] A. Shapiro and T. Homem-de-Mello. A simulation-based approach to two-stage stochastic programming with recourse. *Mathematical Programming*, 81:301–325, 1998.
- [44] A. W. van der Vaart and J.A. Wellner. *Weak Convergence and Empirical Processes*. Springer, 2nd printing 2000 edition, 1996.
- [45] S. van de Geer. *Empirical Processes in M-Estimation*. Cambridge University Press, 2000.
- [46] J. Wang, H. Huang, Z. Luo, and B. Chen. Estimation of covering number in learning theory. In *Proceeding of the Fifth International Conference on Semantics, Knowledge and Grid, 2009*, pages 388–391, Oct 2009.
- [47] Z. Zhang, X. Yang, I. V. Oseledets, G. E. Karniadakis, and L. Daniel. Enabling high-dimensional hierarchical uncertainty quantification by anova and tensor-train decomposition. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 34(1):63–76, Jan 2015.
- [48] D.-X. Zhou. The covering number in learning theory. *J. Complexity*, 18(3):739–767, 2002.