

# Universal $\varepsilon$ -approximators for integrals

Michael Langberg \*

Leonard J. Schulman †

## Abstract

Let  $X$  be a space and  $F$  a family of 0, 1-valued functions on  $X$ . Vapnik and Chervonenkis showed that if  $F$  is “simple” (finite VC dimension), then for every probability measure  $\mu$  on  $X$  and  $\varepsilon > 0$  there is a finite set  $S$  such that for all  $f \in F$ ,  $\sum_{x \in S} f(x)/|S| = [\int f(x)d\mu(x)] \pm \varepsilon$ .

Think of  $S$  as a “universal  $\varepsilon$ -approximator” for integration in  $F$ .  $S$  can actually be obtained w.h.p. just by sampling a few points from  $\mu$ . This is a mainstay of computational learning theory. It was later extended by other authors to families of bounded (e.g.,  $[0, 1]$ -valued) real functions.

In this work we establish similar “universal  $\varepsilon$ -approximators” for families of unbounded nonnegative real functions — in particular, for the families over which one optimizes when performing data classification. (In this case the  $\varepsilon$ -approximation should be multiplicative.)

Specifically, let  $F$  be the family of “ $k$ -median functions” (or  $k$ -means, etc.) on  $\mathbb{R}^d$  with an arbitrary norm  $\varrho$ . That is, any set  $u_1, \dots, u_k \in \mathbb{R}^d$  determines an  $f$  by  $f(x) = (\min_i \varrho(x - u_i))^\alpha$ . (Here  $\alpha \geq 0$ .) Then for every measure  $\mu$  on  $\mathbb{R}^d$  there exists a set  $S$  of cardinality  $\text{poly}(k, d, 1/\varepsilon)$  and a measure  $\nu$  supported on  $S$  such that for every  $f \in F$ ,  $\sum_{x \in S} f(x)\nu(x) \in (1 \pm \varepsilon) \cdot (\int f(x)d\mu(x))$ .

## 1 Introduction

We study numerical integration, the problem of evaluating  $\int f d\mu$ , where

- (a)  $f$  is drawn from a family of nonnegative real-valued functions  $F$  on a metric space  $\mathcal{X} = (X, \varrho)$ ; the only access we need to  $f$  is the ability to evaluate it at points of our choosing.
- (b)  $\mu$  is a probability measure on  $X$ . (In many cases  $\mu$  has finite support of cardinality  $n$ , but this assumption plays no role in our core contributions.)

Our results pertain to families  $F$  of functions, described below, that are of importance in clustering

(classification) and optimization. For these families, we prove a theorem of the following type, where  $\varepsilon$  is any small positive real;  $k$  is a parameter of complexity of the family  $F$ ; and  $\mathcal{X}$  is  $\mathbb{R}^d$  with any norm  $\varrho$ :

For any probability measure  $\mu$  there is a measure  $\nu$  with  $|\text{support}(\nu)| \leq \text{poly}(d, k, 1/\varepsilon)$  such that for all  $f \in F$ ,

$$\int f d\nu \in (1 \pm \varepsilon) \int f d\mu.$$

(Here  $1 \pm \varepsilon$  denotes the interval  $[1 - \varepsilon, 1 + \varepsilon]$ .)

The context for our work lies in strands of research in numerical analysis, statistics and computer science:

(1) *Numerical integration (or quadrature)*. A typical scenario: (a)  $\mathcal{X}$  is Euclidean  $\mathbb{R}^d$ . (b)  $\mu$  is Lebesgue measure restricted to  $[0, 1]^d$ , or to a ball, or  $\mu$  is some other canonical measure such as Gaussian measure about the origin. (c) The family  $F$  is either very restricted (multivariate polynomials, trigonometric functions, etc.) or else subject only to a “tameness” condition (Lipschitz or order of continuity).

Under such conditions, various methods (familiar ones include Newton-Cotes, Gaussian or Clenshaw-Curtis quadrature in one dimension [24], but see also results in high dimension [37]) ensure the existence of measures  $\nu$  with small (carefully chosen) support such that for all  $f \in F$ , either  $\int f d\nu = \int f d\mu$  or  $\int f d\nu \in \int f d\mu \pm \varepsilon$ , depending on whether  $F$  is of the “restricted” or merely “tame” variety.

In this literature the support of  $\nu$  is simply called the set of *evaluation points*.

(2) *Vapnik-Chervonenkis (VC) Theory* in statistics and learning theory. A typical scenario: (a)  $\mathcal{X}$  can be very general but frequently is Euclidean  $\mathbb{R}^d$  or Hamming  $\{0, 1\}^d$ . (b)  $\mu$  is an unknown and arbitrary probability measure on  $X$ ; we (the learners) have access to  $\mu$  only by sampling from it. (c) The range of the functions in  $F$  is  $\{0, 1\}$ , or a bounded-cardinality finite set, or a bounded interval of  $\mathbb{R}$ ; in the first case,  $F$  is simple in the sense that it has bounded *VC Dimension*, while in the latter cases, an appropriate generalization of VC dimension (there are several such) is bounded for  $F$ .

Under such conditions, various theorems [36, 31, 20, 8, 35, 30, 7, 22, 5, 6, 2, 4] (and see [32, 34] for related works) ensure that an “empirical measure”  $\nu$

\*Computer Science Division, Open University of Israel, 108 Ravutski St., Raanana 43107, Israel, email: mikel@openu.ac.il. Work supported in part by The Open University of Israel’s Research Fund (grant no. 46109) and Cisco Collaborative Research Initiative (CCRI).

†California Institute of Technology, MC305-16, Pasadena CA 91125. Email: schulman@caltech.edu. Work supported in part by the NSF.

obtained by selecting a small random sample from  $\mu$  and fixing the uniform measure on those points, has (with high probability) the property that for all  $f \in F$ ,  $\int f d\nu = \int f d\mu \pm \varepsilon$ .

In this literature the support of  $\nu$  is called an  $\varepsilon$ -net,  $\varepsilon$ -transversal,  $\varepsilon$ -sample or  $\varepsilon$ -approximation [8, 27, 9].

**(3) Approximation algorithms for clustering.** A typical scenario: (a)  $\mathcal{X}$  is usually Euclidean (or perhaps  $\ell_1$ )  $\mathbb{R}^d$  but may be also be an explicitly given finite metric. (b)  $\mu$  is usually the uniform measure on an explicitly given finite set of  $n$  points in  $X$ . (c)  $F$  is usually one of the following two families of nonnegative real functions. In each, functions in  $F$  are parameterized by  $k$  points chosen from  $X$ :

$$F_{k\text{-median}} = \{f_{c_1, \dots, c_k}\}_{c_1, \dots, c_k \in X} \text{ where } f_{c_1, \dots, c_k}(x) = \min_{1 \leq i \leq k} \|x - c_i\|$$

$$F_{k\text{-means}} = \{f_{c_1, \dots, c_k}\}_{c_1, \dots, c_k \in X} \text{ where } f_{c_1, \dots, c_k}(x) = \min_{1 \leq i \leq k} \|x - c_i\|^2$$

For these families, various algorithms, some randomized [19, 18, 10] and some deterministic [14], provide a measure  $\nu$  of small support such that  $\int f d\nu \in (1 \pm \varepsilon) \int f d\mu$  for all  $f \in F$ . (This can be a crucial step in clustering, since reduction of the size of the data set allows us to run computationally intensive algorithms.)

In this literature the support of  $\nu$  is usually called a *core-set*, see e.g. [1]. We elaborate on the clustering literature and its relation with our results in a later section of this Introduction.

Since the terminology for  $\nu$  differs between fields, we have adopted  $\varepsilon$ -approximation as the most descriptive of the options. Let  $F$  be a family of nonnegative functions on a metric space  $\mathcal{X}$ .

**DEFINITION 1.** A measure  $\nu$  is an  $\varepsilon$ -approximation for a measure  $\mu$  with respect to  $(F, \mathcal{X})$  if  $\int f d\nu \in (1 \pm \varepsilon) \int f d\mu$  for all  $f \in F$ . For  $g : (0, 1) \rightarrow \mathbb{N}$  nonincreasing,  $(F, \mathcal{X})$  is integrable with complexity  $g$  if for every  $0 < \varepsilon < 1$  and every  $\mu$  there is a  $\nu$  which has  $|\text{support}(\nu)| \leq g(\varepsilon)$  and which is an  $\varepsilon$ -approximation of  $\mu$  w.r.t.  $(F, \mathcal{X})$ . Finally,  $(F, \mathcal{X})$  is finitely integrable if it is integrable with complexity  $g$  for some  $g$ .

This is a multiplicative variant of the uniform Glivenko-Cantelli condition [13, 2].

The question we are interested in is: what classes of functions  $F$  are finitely-integrable? Or, if we consider only measures  $\mu$  of finite support  $n$ , then what classes of functions are integrable with complexity independent of  $n$ ? We can give a partial answer to this question.

**1.1 The general approach: integration by weighted sampling.** Our starting point is the observation (which has been made numerous times previously and falls in the category of “weighted” or “importance”

sampling [11, 23, 3, 33]) that to any probability distribution  $q$  on  $X$  there corresponds the following unbiased estimator for  $\bar{f} = \int f d\mu$ : Sample  $x$  from  $q$ , and set

$$(1.1) \quad T = f(x)\mu(x)/q(x).$$

$T$  is unbiased because  $\int (f(x)\mu(x)/q(x)) dq(x) = \int f(x)\mu(x) dx = \int f(x) d\mu(x)$ . (Technically  $\mathcal{X}$  needs to be finite or compact, and in the latter case,  $q$  in the denominator is a density; nothing depends on these niceties.)

One can, of course, use naïve sampling, i.e.,  $q = \mu$ ; the problem is that the standard deviation of  $T$  can be very large, arbitrarily greater than its expectation. And so the first hurdle we must cross is to reduce  $\text{Var} T$ ; this presents us with a challenging design problem for  $q$ . If we can arrange  $q$  so that  $\text{Var} T$  is not much larger than  $\bar{f}^2$ , then by collecting a small number of samples independently, we can obtain an estimator that is very likely to be within a  $(1 \pm \varepsilon)$  factor of  $\bar{f}$ . This suggests that a plan for constructing an  $\varepsilon$ -approximator  $\nu$  is to sample repeatedly, independently from a carefully chosen  $q$ , then let  $\nu$  be the empirical measure (the uniform distribution on the samples), and integrate  $f\mu/q$  with respect to  $\nu$ .

This brings us, however, to the second and deeper challenge: our estimator needs to *simultaneously* approximate  $\bar{f}$  for each of the infinitely many functions  $f$  in  $F$ . So it is not enough to ensure small probability of error for each  $f$ . This is the “uniformity” challenge that Vapnik and Chervonenkis addressed so successfully in certain (i.e., finite-VC-dimension) families of binary functions. A central part of our work will be to find a substitute argument for the case of multiplicative approximation. Without going into detail we mention that the ideas in the existing additive-approximation extensions of the VC theory do not help with the multiplicative approximation problem. Put simply, that work relies on finitely covering the range of the functions; but  $(0, \infty)$  cannot be covered by finitely many intervals of the form  $(y, (1 + \varepsilon)y)$ .

There are essentially two things we need to do: (1) Quantify the difficulty of integrating a function family  $F$ . We do this using a new parameter which we call the *total sensitivity*  $\mathfrak{S}(F)$ ; this parameter does not have an analogue in the additive approximation theory. (2) Show that an appropriate weighted sampling scheme addresses the “uniformity” challenge, by constructing  $\varepsilon$ -approximations whose size depends on two parameters: one is  $\mathfrak{S}(F)$ , the other is more combinatorial and plays the role analogous to the VC dimension, or more exactly, the shatter function.

**1.2 Our results.** The results of this work are three-fold.

(1) Primarily, we introduce an approach to approximate integration of unbounded nonnegative functions, and show the *existence* of succinct  $\varepsilon$ -approximators for some important families of functions, while also showing that such  $\varepsilon$ -approximators do *not* exist for some other (quite simple) families of functions. To this end we introduce and show the power of the key notion of *sensitivity*; while also showing that it is logically independent of another crucial ingredient, the existence of small  $\varepsilon$ -cover-codes for the family.

We cannot yet provide VC-dimension-style characterizations of when a family of unbounded nonnegative functions is or is not finitely integrable; that is a major open problem and this paper can serve only as a starting point toward its resolution.

(2) We demonstrate the strength of our approach by showing that it yields positive results for a broad family of functions that are important to clustering. What this means is that these families have the off-line/on-line behavior familiar from numerical quadrature: once *someone* has figured out where to put the  $\varepsilon$ -approximator, *anyone* can integrate functions in the family in constant time.

(3) We show that the application in (2) has algorithmic implications: we demonstrate a generic reduction (applicable to any norm on  $\mathbb{R}^d$ ) of the problem of near-optimal clustering of data points, to the problem of bi-criteria approximate clustering. We stress that this connection has already been demonstrated in the literature (especially in [10]); our contribution in this regard is chiefly to show how the idea can be applied generically to any norm and to any clustering exponent  $\alpha$ . (See definition of  $\alpha$  below.)

We now give more detail regarding the families of functions relevant to clustering. These families include and generalize the functions relevant to the well-known  $k$ -means and  $k$ -median problems. Let  $\mathcal{X} = (\mathbb{R}^d, \varrho)$  for an arbitrary norm  $\varrho$ . For  $\alpha > 0$  and  $k \geq 1$ , the “ $k$ -cluster,  $\alpha$ -exponent” function family on  $\mathcal{X}$  is defined as:

**DEFINITION 2.**  $W(\mathcal{X}, k, \alpha) = \{f_{c_1, \dots, c_k}\}_{c_1, \dots, c_k \in \mathbb{R}^d}$  where  $f_{c_1, \dots, c_k}(x) = \min_{1 \leq i \leq k} \varrho(x - c_i)^\alpha$ . (We refer to  $c_1, \dots, c_k$  as the centers of  $f_{c_1, \dots, c_k}$ .)

For example, setting  $\varrho$  to be the Euclidean norm, the  $k$ -median problem on  $\mathcal{X} = (\mathbb{R}^d, \varrho)$  is: Given  $x_1, \dots, x_n \in X$  (and letting  $\mu$  be the uniform measure on  $x_1, \dots, x_n$ ), find  $\Delta(\mathcal{X}, k, 1, \mu) := \inf_{f \in W(\mathcal{X}, k, 1)} \bar{f}$ . Likewise, we obtain the  $k$ -means problem with  $W(\mathcal{X}, k, 2)$ . More generally one seeks for general  $\varrho$  and  $\alpha$ :  $\Delta := \Delta(\mathcal{X}, k, \alpha, \mu) := \inf_{f \in W(\mathcal{X}, k, \alpha)} \bar{f}$ , and frequently one requires as part of the output a spe-

cific  $f^* \in W(\mathcal{X}, k, \alpha)$  for which  $\bar{f}^* = \Delta$ . (It is easy to see that the infimum is achieved.)

As noted above, we show two major results regarding the family  $W(\mathcal{X}, k, \alpha)$ . First, we show that  $W(\mathcal{X}, k, \alpha)$  has succinct  $\varepsilon$ -approximators (which are typically referred to as strong core-sets in the setting of clustering). Second, we describe a general reduction from finding these  $\varepsilon$ -approximators efficiently, to “bi-criteria” approximation of  $W(\mathcal{X}, k, \alpha)$ . Namely, we show for approximation parameters  $c > 0$  and  $\beta > 0$ , that given a function  $f^* \in W(\mathcal{X}, \beta k, \alpha)$  such that  $\bar{f}^* \leq c \min_{f \in W(\mathcal{X}, k, \alpha)} \bar{f}$ , one can efficiently find an  $\varepsilon$ -approximator to  $W(\mathcal{X}, k, \alpha)$ . We note, that using standard ideas, the latter implies efficient (linear in  $n$  time) algorithms for finding near-optimal functions  $f^*$  for the family  $W(\mathcal{X}, k, \alpha)$ . Namely functions  $f^*$  satisfying  $\bar{f}^* \leq (1 + \varepsilon) \min_{f \in W(\mathcal{X}, k, \alpha)} \bar{f}$ . More specifically, one can exhaustively solve the clustering problem on the  $\varepsilon$ -approximator and thus obtain a solution for  $\mathcal{X}$ .

Our results for the function families relevant to clustering can be summarized by the following theorems. In what follows the  $\tilde{O}(\cdot)$  notation neglects logarithmic terms in  $d, k, \alpha$  and  $1/\varepsilon$ .

**THEOREM 1.1.** Let  $\mathcal{X} = (\mathbb{R}^d, \varrho)$ .  $W(\mathcal{X}, k, \alpha)$  has an  $\varepsilon$ -approximator of size  $\tilde{O}\left(\frac{(\alpha^4 + 1)d^2 k^3 2^{4\alpha}}{\varepsilon^2}\right)$ .

In the terms below we assume oracle access to  $f, \varrho$ , the distribution  $\mu$ , and the distribution  $q$  we construct.

**THEOREM 1.2.** Let  $X \subset \mathbb{R}^d$  be of size  $n$ . Let  $\mathcal{X} = (X, \varrho)$ . Let  $\beta > 0$  and  $c > 0$  be constants. Let  $\mathcal{A}$  be any algorithm that finds a function  $f^* \in W(\mathcal{X}, \beta k, \alpha)$  such that  $\bar{f}^* \leq c \min_{f \in W(\mathcal{X}, k, \alpha)} \bar{f}$  in time  $T_{\beta, c}$ . Then an  $\varepsilon$ -approximator of size  $\tilde{O}\left(\frac{(\alpha^4 + 1)d^2 k(c + \beta k)^2 2^{6\alpha}}{\varepsilon^2}\right)$  can be found in time  $T_{\beta, c} + O(ndk\beta) + \tilde{O}\left(\frac{(\alpha^4 + 1)d^2 k(c + \beta k)^2 2^{6\alpha}}{\varepsilon^2}\right)$ .

**1.3 Related work.** As mentioned above, strong core-sets (or  $\varepsilon$ -approximators) for  $k$ -means and  $k$ -median have been extensively studied in the past. In [19], it is shown that for sets  $X \subset \mathbb{R}^d$  of size  $n$ , a core-set of size  $O(k\varepsilon^{-d} \log n)$  can be constructed. The dependence on  $n$  was removed in [18] to obtain core-sets of size  $O(k^3 \varepsilon^{-d-1})$ . For high dimensional spaces, Chen [10] reduced the exponential dependence on  $d$  at the price of a logarithmic dependence in  $n$  and obtains core-sets of size  $O(k^2 d \varepsilon^{-2} \log n)$ . For comparison, the core-sets we obtain are of size polynomial in  $d, k$ , and  $1/\varepsilon$ , and independent of  $n$ . (It is possible such a result may be obtained in other ways; Feldman (private communication) has noted that a variation in the work of [15] may yield such core-sets.)

In a nutshell, the results above all use the paradigm of “bi-criteria” approximation mentioned previously. Namely, they show how to construct strong core-sets given a solution to the  $\beta k$ -median (or mean) of value comparable to the optimal  $k$  clustering. This then leads to efficient algorithms for the solution of  $k$ -means and  $k$ -median. The best running time to date is that of  $\tilde{O}(ndk + 2^{\tilde{O}(k/\varepsilon)} + d \text{poly}(k/\varepsilon))$  presented in [15].

In its overall structure, our algorithms for the family  $W(\mathcal{X}, k, \alpha)$  closely resemble those mentioned above. Namely, given a bi-criteria approximation, a certain random process is preformed in order to obtain a small  $\varepsilon$ -approximation.

For comparison of our work with earlier works mentioned previously in VC theory establishing small  $\varepsilon$ -approximations, note that all these works provide what is essentially an additive  $\varepsilon$ -approximation. The requirement of multiplicative approximation is what makes our work very different from those referenced earlier. (However, there are interesting connections to the VC theory. The dependence on  $d$  in our argument relies upon bounding a shatter-like function.)

Finally we note that (multiplicative)  $\varepsilon$ -approximators for finite sets  $X$  and for a number of function families  $F$  defined on the line were studied in [17]. These families  $F$  include linear functions, piecewise monotone functions, and a variant of the  $k$ -median problem on the line. Most of these function families do not have finite  $\varepsilon$ -approximators when  $X$  is taken to be infinite. The work of [17] focuses on quantifying the dependence of the  $\varepsilon$ -approximator’s size on the cardinality  $n$  of  $X$ .

**1.4 Layout of the paper.** In Section 2 we introduce the notion of total sensitivity and show how it can be used to design our sample distribution  $q$ . We then focus on the families  $F = W(\mathcal{X}, k, \alpha)$ : In Section 3 we analyze their total sensitivity, and in Sections 4 and 5 we use this analysis to prove Theorems 1.1 and 1.2. Several of our proofs are omitted and can be found in a full version of the paper [26].

## 2 Total sensitivity

**2.1 Sensitivity and sampling.** For each  $x \in X$ , define the *sensitivity* of  $x$  w.r.t.  $(F, \mu)$  by  $\sigma_{F, \mu}(x) = \sup_{f \in F} f(x)/\bar{f}$ ; when the identities of  $F$  and  $\mu$  are either clear or immaterial we abbreviate this by  $\sigma(x)$ . Define the *total sensitivity* of  $F$  by  $\mathfrak{S}(F) = \sup_{\mu} \int \sigma_{F, \mu}(x) d\mu(x)$ . From a theoretical perspective these are the key concepts but it will be useful at a later point if the results of the current section are stated in terms of the following quantities:  $s_{F, \mu}(x)$  (or  $s(x)$ ) is any upper bound on  $\sigma_{F, \mu}(x)$  (or  $\sigma(x)$ ), and  $S(F) =$

$\sup_{\mu} \int s_{F, \mu}(x) d\mu(x)$ . (Obviously  $S(F) \geq \mathfrak{S}(F)$ .) The inequality to keep in mind is: for any  $f \in F$ , any  $\mu$  and any  $x \in X$ ,

$$(2.2) \quad f(x) \leq s(x)\bar{f}.$$

The beauty of the notion of sensitivity is that it leads to a judicious weighted sampling scheme for integration. Let

$$(2.3) \quad q(x) = s(x)\mu(x)/S(F)$$

and let  $T$  be the estimator for  $\bar{f}$  given in Equation 1.1. Observe that the sampling process tends to pick  $x$ ’s of high sensitivity. Now we’ll see why total sensitivity is a crucial parameter:

**THEOREM 2.1.**  $\text{Var } T \leq (S(F) - 1)\bar{f}^2$ .

*Proof.* Let  $S := S(F)$ . We analyze the value of  $(1/\bar{f}^2) \text{Var } T$  which is equal to

$$\begin{aligned} & (1/\bar{f}^2) \int \left( \frac{f(x)\mu(x)}{q(x)} - \bar{f} \right)^2 dq(x) \\ &= (1/\bar{f}^2) \int \frac{\mu(x)s(x)}{S} \left( \frac{f(x)S}{s(x)} - \bar{f} \right)^2 dx \\ &= (1/\bar{f}^2) \int \left( \frac{f(x)^2\mu(x)S}{s(x)} - 2\bar{f}f(x)\mu(x) + \frac{\mu(x)s(x)\bar{f}^2}{S} \right) dx \\ &= \left( (1/\bar{f}^2) \int \frac{f(x)^2S}{s(x)} d\mu(x) \right) - \left( (2/\bar{f}) \int f(x) d\mu(x) \right) \\ & \quad + \left( \int s(x)/S d\mu(x) \right) \\ &= (1/\bar{f}^2) \int \frac{f(x)^2S}{s(x)} d\mu(x) - 2 + 1 \end{aligned}$$

Now we apply Inequality 2.2 to obtain,

$$\leq (1/\bar{f}^2) \int f(x)\bar{f}S d\mu(x) - 1 = (S/\bar{f}) \int f(x) d\mu(x) - 1 = S - 1.$$

The effect of averaging over multiple samples is expressed by a simple application of Chebyshev’s inequality.

**LEMMA 2.1.** *Let  $\varepsilon > 0$ . Let  $f \in F$ . Let  $R$  be a random sample of  $X$  of size  $a \geq \frac{2(S-1)}{\varepsilon^2}$  according to the distribution  $q$ . Then  $\Pr \left[ \left| \bar{f} - \frac{1}{a} \sum_{x \in R} \frac{S(F)f(x)}{s(x)} \right| \geq \varepsilon \bar{f} \right] \leq 1/2$ .*

**2.2 Case studies.** To put our definitions in perspective we present some examples.

**2.2.1 Families that do not have  $\varepsilon$ -approximators.** In order to illustrate the nature of the problem we examine a few very simple families  $F$  that do not have finite  $\varepsilon$ -approximators.

The simplest example is this:  $X$  is the unit interval and  $F_{\text{ivl}} := \{f_{a,b}\}_{a < b}$  where  $f_{a,b}(x) = 1$  if  $a < x < b$  and 0 otherwise. This is a family that trivially has an  $\varepsilon$ -approximator for additive approximation: let  $\nu$  be supported uniformly on the multiset  $x_k = \inf\{x : \int_0^x \mu(x) dx \geq k\varepsilon/2\}$ . The conclusion can also be seen from the more abstract consideration that  $F_{\text{ivl}}$  has VC dimension 2.

On the other hand, if  $\mu$  is any measure without atoms, then for any finitely supported  $\nu$ , there is a pair  $a < b$  such that  $\mu((a, b)) > 0$  but  $(a, b) \cap \text{support } \nu = \emptyset$ , so  $(\int f_{a,b} d\nu)/(\int f_{a,b} d\mu) = 0$ . It is easy to see that  $\mathfrak{S}(F_{\text{ivl}}) = \infty$  because every point  $x$  has unbounded sensitivity.

For our second example we again let  $X$  be the unit interval and let  $F_{\text{ray}} := \{f_a\}_a$  where  $f_a(x) = 1$  if  $a < x$  and 0 otherwise. For the purpose of additive approximation this is an even easier family than  $F_{\text{ivl}}$  (it has VC dimension 1), but this family too has no finite  $\varepsilon$ -approximators: again let  $\mu$  be the uniform measure. If  $\nu$  is finitely supported, let  $a$  be the greatest point in its support; then  $(\int f_a d\nu)/(\int f_a d\mu) = 0$ . To see that  $\mathfrak{S}(F_{\text{ray}}) = \infty$ , let  $\mu$  be the uniform measure. Now  $\sigma(x) \geq \sup_{a < x} 1/(1-a) = 1/(1-x)$ , so  $\mathfrak{S}(F_{\text{ray}}) \geq \int_0^1 1/(1-x) dx = \infty$ .

It is worth examining the “finitary” versions of each of these examples, with  $\mu$  required to have support of cardinality  $\leq n$ . For  $F_{\text{ivl}}$ , simply take  $\mu$  to be uniform on  $n$  points  $x_1 < \dots < x_n$ . None of these points may be omitted in an  $\varepsilon$ -approximator. And,  $\mathfrak{S}(F_{\text{ivl}})$  (under the restriction to measures of support  $\leq n$ ) equals  $n$ . For  $F_{\text{ray}}$ , take the same points  $x_1 < \dots < x_n$  but now set  $\mu(x_i) = 2^{-i}$  (except that  $\mu(x_1) = 2^{-1} + 2^{-n}$ ). Then  $\sigma(x_i) \geq 2^{i-1}$ , so  $\mathfrak{S}(F_{\text{ray}})$  (again under the cardinality restriction) is  $\geq \sum_i \sigma(x_i)\mu(x_i) \geq \sum_i 2^{i-1}2^{-i} = n/2$ .

**2.2.2 Families that have  $\varepsilon$ -approximators, yet sampling from  $\mu$  is futile.** Our next example is the simplest and most tractable prototype of the families of functions that we shall show in this paper to have finite  $\varepsilon$ -approximators. Let  $X = \mathbb{R}$  and let  $F_{1\text{-means}} := \{f_a\}_{a \in \mathbb{R}}$  where  $f_a(x) = (x-a)^2$ . Let’s see why trying to construct an  $\varepsilon$ -approximator by sampling repeatedly from  $\mu$  and setting  $\nu$  to be the empirical measure, fails.

For  $0 < p < 1$ , consider the following measure  $\mu$  supported on  $\{0, 1\}$ :  $\mu(\{0\}) = p, \mu(\{1\}) = 1-p$ . Observe that  $\int f_1 d\mu = p$  but that a sample of  $\sim 1/p$  points is needed for  $\nu$  to be likely to include 0 in its support; so long as this does not occur,  $\int f_1 d\nu = 0$ . So there is no finite sample size at which we likely obtain an  $\varepsilon$ -approximator for every  $\mu$ .

However, it is plain that a weighted sampling scheme  $q$  (as described in Sec. 1.1) can be made to

work on these simple “counterexample” measures, for instance by letting  $q$  be uniform on  $\{0, 1\}$ . This leaves open the question of handling general measures  $\mu$ . In Section 2.4 we show that indeed weighted sampling can generate  $\varepsilon$ -approximators for the “1-means” family for any  $\mu$  and in any dimension, and we will actually deduce the sampling scheme  $q$  in closed form.

For the more complicated (and more important) families of functions that are important for clustering, there is almost certainly no closed form for optimal solutions. Nonetheless we show in Section 4 that it is possible to obtain an effective sampling scheme.

**2.3 Properties of total sensitivity.** What operations on function families preserve bounds on total sensitivity? Two are easy to see: addition and projective closure.

**PROPOSITION 1.** *Given families of nonnegative functions  $F, G$  on  $\mathcal{X}$ , let  $F + G = \{f + g : f \in F, g \in G\}$ . Then  $\mathfrak{S}(F + G) \leq \mathfrak{S}(F) + \mathfrak{S}(G)$ .*

*Proof.*  $(f(x) + g(x))/(\bar{f} + \bar{g}) \leq f(x)/\bar{f} + g(x)/\bar{g}$  (keep in mind positivity of the quantities).

Given a family of nonnegative functions  $F$  on  $\mathcal{X}$ , its *closure*  $F^c$  includes any function  $g$  such that for any bounded set  $A \subseteq X$  and any  $\delta > 0$  there exist  $f \in F$  and  $c > 0$  such that for all  $x \in A$ ,  $cf(x) \leq g(x) \leq e^\delta cf(x)$ . The *projective closure* of  $F$ ,  $\mathbb{P}F^c$ , is formed by choosing one representative from each equivalence class in  $F^c$  w.r.t. multiplication by positive scalars. (We need this operation mainly in order to automatically include constant functions in our families.) The proof of the following proposition is omitted.

**PROPOSITION 2.**  $\mathfrak{S}(\mathbb{P}F^c) = \mathfrak{S}(F)$ .

Finally, multiplication does *not* preserve good bounds on total sensitivity. Let  $FG = \{fg : f \in F, g \in G\}$ . The proof of the following proposition is omitted.

**PROPOSITION 3.** *There exist  $F, G$  for which  $\mathfrak{S}(FG) \geq e^{\Omega(\max\{\mathfrak{S}(F), \mathfrak{S}(G)\})}$ .*

**2.4 Vector norm functions.** The key to calculating the total sensitivity of  $F_{1\text{-means}}$  (mentioned above) is the observation that there is a real vector space  $V$  (the vector space of affine linear functions on  $\mathbb{R}^1$ ) such that for every  $f \in F_{1\text{-means}}$  there is a  $v \in V$  satisfying  $f(x) = |v(x)|^2$ .

**THEOREM 2.2.** *Let  $V \subseteq \mathbb{R}^X$  be a real vector space of dimension  $d$ . Let  $F = \{|p(x)|^2 : p \in V\}$ . Then  $\mathfrak{S}(F) = d$ .*

*Proof.* Let  $\mu$  be any probability measure on  $X$ , and consider it as defining an inner product on  $V$  that is diagonal in the “basis” of delta-functions: the inner product of  $u, v \in V$  is  $\int u(x)v(x) d\mu(x)$ . (Correspondingly, we freely consider  $\mu$  as a diagonal matrix.) Let  $p^1, \dots, p^d$  be an orthonormal basis for  $V$  w.r.t.  $\mu$ , and let  $P$  be the (possibly infinite) matrix having these vectors as rows. Consider  $\mathbb{R}^d$  as a space of column vectors and let  $S^{d-1}$  denote the unit sphere (w.r.t. the identity inner product). Observe that

$$\sigma(x) = \sup_{u \in S^{d-1}} \frac{|(u^\dagger P)(x)|^2}{u^\dagger P \mu P^\dagger u} = \sup_{u \in S^{d-1}} |(u^\dagger P)(x)|^2$$

The last term is maximized by letting  $u$  be a scalar multiple of the column vector  $P(x) = (p^1(x), \dots, p^d(x))$ . So,

$$\sigma(x) = \left| \frac{P^\dagger(x)P(x)}{\|P(x)\|} \right|^2 = \|P(x)\|^2$$

Therefore  $\mathfrak{S} = \int \sigma d\mu = \text{Tr } P^\dagger P \mu = \text{Tr } P \mu P^\dagger$ . Since the vectors  $p^i$  are orthonormal w.r.t.  $\mu$ ,  $\mathfrak{S} = d$ .

**COROLLARY 2.1.** *The family of squares of  $d$ -variate real polynomials of total degree  $\leq t$  has total sensitivity  $\binom{t+d}{d}$ .*

Observe that the vector space of affine linear functions on  $\mathbb{R}$ ,  $V_{\text{aff}} = \{b + ax\}_{a,b \in \mathbb{R}}$ , is of dimension 2. Let  $F_{d,1}$ -means be  $\{f_a\}_{a \in \mathbb{R}^d} : \mathbb{R}^d \rightarrow \mathbb{R}$  where  $f_a(x) = \sum (x_i - a_i)^2$ . We can see that  $\mathfrak{S}(F_{d,1}\text{-means}) \leq 2d$  by applying Corollary 2.1 and Proposition 1; however, we shortly improve this bound.

Corollary 2.1 and Proposition 1 together beg the question whether one can bound the total sensitivity of the cone of nonnegative real polynomials of total degree  $\leq 2t$ . Proposition 1 shows that for the cone of “sum of squares” polynomials, the bound  $\binom{t+d}{d}k$  holds for polynomials which are sums of  $k$  squares; we do not know whether the dependence on  $k$  is necessary. Moreover, there are nonnegative polynomials which are not sums of squares (except in the special cases  $d = 1$ ;  $t = 1$ ; and  $(d = 2, t = 2)$ ) [21]).

### 3 Clustering functions and their total sensitivity

In this section we study the total sensitivity of the family  $F = W(\mathcal{X}, k, \alpha)$ . Let  $\mathcal{X} = (X, \rho)$ . Consider any probability measure  $\mu$  on  $X$ ; let  $\bar{f} = \int f d\mu$  and  $\Delta = \inf_{f \in F} \bar{f}$ .

**THEOREM 3.1.** (1) For  $\alpha \geq 1$ ,  $\mathfrak{S}(W(\mathcal{X}, k, \alpha)) \leq (k + 1)2^{2\alpha} + 2^\alpha$ . (b) For  $0 \leq \alpha \leq 1$ ,  $\mathfrak{S}(W(\mathcal{X}, k, \alpha)) \leq 2k + 3 + 2\sqrt{6k}$ .

Observe that these bounds are independent of  $\mathcal{X}$ .

*Proof.* If  $\alpha = 0$  the theorem is trivial. Otherwise, let  $f^* = f_{u_1^*, \dots, u_k^*}$  be a function in  $W(\mathcal{X}, k, \alpha)$  for which  $\bar{f}^* = \Delta$ . Let  $U_i$  be the Voronoi cell of  $u_i^*$ , and let  $p_i = \mu(U_i)$ . (Each  $p_i$  is positive unless the support of  $\mu$  has cardinality less than  $k$  in which case the theorem is trivial.) Let  $m_i = \frac{1}{p_i} \int_{U_i} \varrho(x - u_i^*)^\alpha d\mu(x)$ , so  $\Delta = \sum p_i m_i$ . By a simple Markov inequality, for each  $i$ ,  $\mu(B(u_i^*, (2m_i)^{1/\alpha}) \cap U_i) \geq p_i/2$ . (Here  $B(x, r)$  denotes the closed ball of radius  $r$  about  $x$ .)

We now analyze  $\mathfrak{S}$ . Let  $x \in U_i$ , and let  $f = f_{u_1, \dots, u_k}$  be any function in  $W(\mathcal{X}, k, \alpha)$ . Let  $u$  denote a closest point to  $u_i^*$  in  $\{u_1, \dots, u_k\}$ , and let  $\varrho_i = \varrho(u - u_i^*)$ . Then  $\bar{f} \geq \int_{U_i} f d\mu \geq [\max(0, \varrho_i - (2m_i)^{1/\alpha})]^\alpha \frac{p_i}{2}$ . Also,  $\bar{f} \geq \Delta$ . Thus, for a parameter  $q \in [0, 1]$ ,  $\bar{f} \geq [\max(0, \varrho_i - (2m_i)^{1/\alpha})]^\alpha \frac{qp_i}{2} + (1 - q)\Delta$ . At this point the arguments for Parts (1,2) of the theorem diverge. Part (1):

$$\begin{aligned} \sigma(x) &= \max_f f(x)/\bar{f} \leq \max_f \varrho(x - u)^\alpha / \bar{f} \\ &\leq \max_f \frac{(\varrho_i + \varrho(x - u_i^*))^\alpha}{[\max(0, \varrho_i - (2m_i)^{1/\alpha})]^\alpha qp_i/2 + (1 - q)\Delta} \\ &\leq \max_{\varrho_i \geq 0} \frac{2^{\alpha-1} (\varrho_i^\alpha + \varrho(x - u_i^*)^\alpha)}{[\max(0, \varrho_i - (2m_i)^{1/\alpha})]^\alpha qp_i/2 + (1 - q)\Delta} \\ &\quad (\text{apply } |a + b|^\alpha \leq 2^{\alpha-1}(|a|^\alpha + |b|^\alpha)) \\ &= \max_{\varrho_i \geq (2m_i)^{1/\alpha}} \frac{2^{\alpha-1} (\varrho_i^\alpha + \varrho(x - u_i^*)^\alpha)}{(\varrho_i - (2m_i)^{1/\alpha})^\alpha qp_i/2 + (1 - q)\Delta} \\ &\leq \max_{\varrho_i^\alpha \geq 2m_i} \frac{2^{\alpha-1} (\varrho_i^\alpha + \varrho(x - u_i^*)^\alpha)}{\max\{0, 2^{1-\alpha} \varrho_i^\alpha - 2m_i\} qp_i/2 + (1 - q)\Delta} \\ &\quad (\text{apply } |a - b|^\alpha \geq \max\{0, 2^{1-\alpha}|a|^\alpha - |b|^\alpha\}) \end{aligned}$$

Let  $G(\varrho_i^\alpha) = \frac{\varrho_i^\alpha + \varrho(x - u_i^*)^\alpha}{(2^{1-\alpha} \varrho_i^\alpha - 2m_i) qp_i/2 + (1 - q)\Delta}$ . Observe that  $\text{sign}\left(\frac{\partial G}{\partial \varrho_i^\alpha}\right)$  is independent of  $\varrho_i^\alpha$  and thus  $G$  is monotone as a function of  $\varrho_i^\alpha$ . Moreover,  $\frac{\varrho_i^\alpha + \varrho(x - u_i^*)^\alpha}{(1 - q)\Delta}$  is increasing as a function of  $\varrho_i^\alpha$ . Thus the bound on  $\sigma(x)$  is maximized at either  $\varrho_i = 2m_i^{1/\alpha}$  or  $\varrho_i = \infty$ . We conclude that for  $x \in U_i$ :

$$\begin{aligned} \sigma(x) &\leq \max\left(\frac{2^{2\alpha-1} m_i + 2^{\alpha-1} \varrho(x - u_i^*)^\alpha}{(1 - q)\Delta}, \frac{2^{2\alpha-1}}{qp_i}\right) \\ &\leq \frac{2^{2\alpha-1} m_i + 2^{\alpha-1} \varrho(x - u_i^*)^\alpha}{(1 - q)\Delta} + \frac{2^{2\alpha-1}}{qp_i} \end{aligned}$$

Thus  $\int \sigma d\mu$  is equal to

$$\begin{aligned} &\sum_i \left( \int_{U_i} \sigma d\mu \right) \\ &\leq \sum_i \left( \int_{U_i} \frac{2^{2\alpha-1} m_i + 2^{\alpha-1} \varrho(x - u_i^*)^\alpha}{(1 - q)\Delta} + \frac{2^{2\alpha-1}}{qp_i} d\mu \right) \end{aligned}$$

$$\begin{aligned}
&= \sum_i \left( \frac{2^{2\alpha-1} p_i m_i}{(1-q)\Delta} + \frac{2^{\alpha-1} p_i m_i}{(1-q)\Delta} + \frac{2^{2\alpha-1}}{q} \right) \\
&= \frac{2^{2\alpha-1} + 2^{\alpha-1}}{1-q} + \frac{2^{2\alpha-1} k}{q} \\
&\leq \left( \sqrt{2^{2\alpha-1} + 2^{\alpha-1}} + \sqrt{2^{2\alpha-1} k} \right)^2 \\
&\quad \text{(the minimum of } a/(1-q) + b/q \text{ is } (\sqrt{a} + \sqrt{b})^2) \\
&\leq (k+1)2^{2\alpha} + 2^\alpha \quad \text{(use } (a+b)^2 \leq 2(a^2 + b^2))
\end{aligned}$$

This suffices to bound  $\mathfrak{S} = \sup_\mu \int \sigma(x) d\mu(x)$  and to prove Part (1) of the theorem. Part (2) of the proof is similar in nature and is omitted.

#### 4 $\varepsilon$ -approximators, arrangements and covering codes for $W(\mathcal{X}, k, \alpha)$

Let  $\mathcal{X} = (\mathbb{R}^d, \varrho)$ . In this section we study the family  $F = W(\mathcal{X}, k, \alpha)$  and show that it has succinct  $\varepsilon$ -approximators. Let  $S$  be a bound on the total sensitivity of  $F$ . We prove the following theorem which is a rephrased version of Theorem 1.1 stated in the Introduction.

**THEOREM 4.1.** *Let  $\mathcal{X} = (\mathbb{R}^d, \varrho)$ . The family  $W(\mathcal{X}, k, \alpha)$  has an  $\varepsilon$ -approximator of size  $O\left(\frac{(\alpha^4+1)d^2 k S^2}{\varepsilon^2} \log\left(\frac{(\alpha+1)dkS}{\varepsilon}\right)\right)$ . Thus by Theorem 3.1, for  $\alpha \geq 1$  the  $\varepsilon$ -approximator is of size  $\tilde{O}\left(\frac{\alpha^4 d^2 k^3 2^{4\alpha}}{\varepsilon^2}\right)$  and for  $\alpha \leq 1$  the  $\varepsilon$ -approximator is of size  $\tilde{O}\left(\frac{d^2 k^3}{\varepsilon^2}\right)$ . Here  $\tilde{O}(\cdot)$  neglects logarithmic terms in  $d, k, \alpha$  and  $1/\varepsilon$ .*

Roughly, speaking this is done in three major steps. Primarily, recall that  $\mathcal{X}$  (and thus  $F$ ) is defined by a norm  $\varrho$  which in turn is defined by a centrally symmetric convex set  $C_\varrho$ . In our first step, we show that one may assume w.l.o.g. that  $\varrho$  is *well behaved*. Namely, that any arrangement (in the sense of combinatorial geometry) described by  $n$  translates of  $C_\varrho$  (i.e., the dissection of  $\mathbb{R}^d$  by the boundaries of the translates) has low complexity, of approximately  $n^d$ . In our second step, we define the notion of an “ $\varepsilon$ -cover code” for the family  $F$ . Namely, a set of functions  $F' \subset F$  that *approximates* the set  $F$  with respect to any finite subset of the support  $X$ . We show that  $F$  has small cover-codes  $F'$  if its underlying norm  $\varrho$  is well behaved. Finally, in our third step we show that a small  $\varepsilon$ -cover code for  $F$  implies a small  $\varepsilon$ -approximator. This will conclude the proof of Theorem 4.1. We give a detailed outline below. The detailed proof of several of the following theorems is omitted.

##### 4.1 Well behaved norms.

**DEFINITION 3.** *Let  $\varrho$  be a norm corresponding to the centrally symmetric convex set  $C$ . Consider any collec-*

*tion of convex sets  $\{C_1, \dots, C_n\}$  where each set  $C_i$  is equal to  $r_i C_\varrho + v_i$ . Here  $v_i$  is a vector in  $\mathbb{R}^d$ ,  $r_i$  is a positive real, and  $r_i C = \{x \mid \varrho(x) \leq r_i\}$ . The collection of sets  $C_i$  describes an arrangement in  $\mathbb{R}^d$ . We say that  $\varrho$  is  $\Gamma$ -well behaved if the complexity of any such arrangement is bounded by  $(n\Gamma)^d$ .*

**THEOREM 4.2.** *Let  $\varrho$  be any norm. Let  $\Gamma = (cd/\sqrt{\varepsilon})^d$  for a sufficiently large constant  $c$ . There exists a  $\Gamma$ -well behaved norm  $\varrho'$  such that any  $\varepsilon$ -approximator for  $\varrho'$  will yield an  $O((1+\varepsilon)^\alpha - 1)$ -approximator for the original  $\varrho$ .*

Notice that when  $\varepsilon$  is small compared to  $1/\alpha$ , the quality loss in our approximator is small.

In a nutshell, the proof of Theorem 4.2 is based on a previous work of ours [25] in which certain approximations to convex sets  $C$  were studied. Namely, using a proof technique similar to that of Dudley for convex shape approximation by a polytope with few vertices [12], in Theorem 1 of [25] it is shown that every centrally symmetric convex set  $C$  has a corresponding centrally symmetric convex set  $C' \subseteq C$  such that (a)  $C'$  is a polyhedron of low complexity, and (b) for any  $x \in \mathbb{R}^d$  it holds that  $(1-\varepsilon)\|x\|_{C'} \leq \|x\|_C \leq \|x\|_{C'}$ . Here,  $\|x\|_C = \inf\{r > 0 \mid \frac{x}{r} \in C\}$ . More specifically, it was shown in [25] that  $C'$  has at most  $(cd/\sqrt{\varepsilon})^d (d-1)$ -dimensional facets, where  $c$  is a sufficiently large universal constant. Defining  $\varrho'$  corresponding to  $C'$  and using a few additional ideas suffices to prove Theorem 4.2.

(This is the place to draw attention to a crucial distinction between Euclidean norm and general norms  $\varrho$ . Recall that the collection of positive homothetes (= translation and multiplication by positive scalars) of the unit ball in  $\mathbb{R}^d$ , has VC dimension  $d+1$ . This implies that the complexity of an arrangement of  $n$  balls is  $O(n^{d+1})$ . Grünbaum [16] conjectured that the same VC dimension bound held for the positive homothetes of any fixed compact convex set in  $\mathbb{R}^d$ . If this (or even a weaker  $O(d)$  upper bound) were true, the size of our  $\varepsilon$ -approximators in Theorem 4.1 would no longer have quadratic dependence in  $d$  but rather linear dependence. However, the conjecture has been falsified. The first counterexample was due to Naiman and Wynn, who showed that the collection of translates of a box [28] has VC dimension  $\lfloor 3d/2 \rfloor$ . This of course is not large enough to be an obstacle to our application. But very recently, Naszódi and Taschuk showed that in dimension 3 and above there is *no* upper bound on the VC dimension of the positive homothetes of convex bodies [29]. Actually even stronger, there is a convex body whose collections of translates has infinite VC dimension.)

**4.2  $\varepsilon$ -cover codes.** A central role in the VC theory of additive approximation is played by the *shatter function* of a family  $F$  of Boolean functions. (Frequently, VC dimension appears merely as a means to bound the shatter function.) The viewpoint which generalizes to our situation is this: Shatter functions measure the size of *covering codes* or *transversals* for (restrictions of)  $F$ . In the boolean case the notion of “covering” is trivial: one function covers another only if their restrictions are identical. But we require something more general. Let  $F$  be a function family and  $s(x)$  a sensitivity bound (as in Section 2.1). Let  $A \subseteq X$  be finite,  $a = |A|$ . For  $g : X \rightarrow \mathbb{R}$ , let  $\nu_A(g) = (1/a) \sum_{x \in A} g(x)$ . For  $f, f' \in F$  and  $x \in A$  define  $\hat{f} = \nu_A(f/s)$  and

$$D_{A,x}(f, f') = \left| \frac{f(x)}{\hat{f}s(x)} - \frac{f'(x)}{\hat{f}'s(x)} \right|$$

Notice that  $D_{A,x}(f, f')$  also depends on our bound  $s(x)$  on the sensitivity of  $F$  at  $x$ , however, we do not write  $s$  explicitly as a parameter in  $D$ . The definition of  $D$  and that which follows are designed to fit our needs in Theorem 4.4 of Section 4.3.

**DEFINITION 4.**  $F' \subseteq F$  is an  $\varepsilon$ -cover-code for  $(F, A, s)$  if for every  $f \in F$  there is an  $f' \in F'$  such that  $\frac{\bar{f}'}{\hat{f}'} \leq \frac{\bar{f}}{\hat{f}}$  and for every  $x \in A$ ,  $D_{A,x}(f, f') \leq \frac{\varepsilon}{64S} \left( 1 + \frac{f(x)}{\hat{f}s(x)} + \frac{f'(x)}{\hat{f}'s(x)} \right)$ .

**THEOREM 4.3.** Let  $\mathcal{X} = (\mathbb{R}^d, \varrho)$ . Let  $F = W(\mathcal{X}, k, \alpha)$  where  $\varrho$  is a  $\Gamma$ -well behaved norm. Let  $A \subset \mathbb{R}^d$  be a set of size  $a$ ,  $F$  has an  $\varepsilon$ -cover-code  $F'$  of size  $\left[ \left( \frac{Sa}{\varepsilon} \right)^{\Theta(\alpha^2+1)} \Gamma \right]^{2dk}$ .

Roughly speaking, our proof of Theorem 4.3 has two steps. In the first step we show the existence of a small set of functions  $G$  for which for any  $f \in F$  there exists a constant  $c_f$  and a function  $g \in G$  which covers  $f$ . More specifically, for any  $x \in A$ :

$$(4.4) \quad \left| \frac{f(x)}{s(x)} - \frac{c_f g(x)}{s(x)} \right| \leq \frac{\varepsilon}{256S} \hat{f}.$$

This step is rather technically involved, and is based on a certain decomposition of the set  $A$  in  $\mathcal{X} = (\mathbb{R}^d, \varrho)$ .

We do not take  $G$  to be our  $\varepsilon$ -cover-code since  $G$  is not a subfamily of  $F$ . In our second step we define a mapping from  $G$  into  $F$ . Namely, we partition the set  $F$  into  $|G|$  disjoint sets, with  $F^g$  being the set of functions  $f$  covered by  $g$ . We then set  $f^g \in F^g$  to be a function in  $F^g$  minimizing  $\bar{f}/\hat{f}$ . Namely,  $\bar{f}^g/\hat{f}^g \leq \bar{f}/\hat{f}$  for all  $f \in F^g$ . We show that any  $f$  covered by  $g$  will also be covered by  $f^g$ . The set  $F' = \{f^g \mid g \in G\}$

is our final cover code. The proof of this second step follows (almost) directly from the properties of  $G$  stated in Equation 4.4.

**4.3  $\varepsilon$ -approximators.** Finally, we turn to show a general VC-type argument. Loosely speaking, we show that small cover codes imply succinct  $\varepsilon$ -approximators via random sampling. The analysis of our theorem holds for any family  $F$  with support  $X$  and total sensitivity at most  $S$ .

**THEOREM 4.4.** Suppose that for some  $a \geq 8(S - 1)/\varepsilon^2$ , every  $A \subseteq X$  of cardinality  $|A| = 2a$  possesses an  $\varepsilon$ -cover-code (w.r.t.  $F$  and sensitivity bound  $s$ ) of cardinality at most  $\frac{1}{8} e^{\frac{a\varepsilon^2}{100S^2}}$ . Then a sample of  $a$  points from  $q$  is (with probability  $\geq 1/2$ ) an  $\varepsilon$ -approximator for  $F$ .

The proof of Theorem 4.4 has the following outline. Let  $E_0$  be the event that a sample  $R$  of size  $a$  according to the distribution  $q$  is not an  $\varepsilon$ -approximator for  $F$ . Say that  $(f, R)$  is  $\varepsilon$ -bad if  $|\bar{f} - \nu_R(Sf/s)| > \varepsilon f$ . (Recall, for a function  $g$ , that  $\nu_R(g) = (1/|R|) \sum_{x \in R} g(x)$ .) So,  $E_0$  occurs if there is an  $f \in F$  s.t.  $(f, R)$  is  $\varepsilon$ -bad. Our objective is to bound  $\Pr[E_0]$ .

In the manner of Vapnik and Chervonenkis, now let  $G$  be an additional multiset of size  $a$  chosen independently at random according to  $q$ . Let  $E_1$  be the event “ $\exists f \in F : (f, R)$  is  $\varepsilon$ -bad and  $(f, G)$  is  $\varepsilon/2$ -good” (here good is the complement of bad). Using Lemma 2.1, it is not hard to verify that  $\Pr[E_1] \leq \Pr[E_0] \leq 2 \Pr[E_1]$ . And thus we turn to bound  $\Pr[E_1]$ .

Let  $A$  be an independent random sample of size  $2a$  according to the distribution  $q$ . Let  $R$  be a random (uniform) sample from  $A$  of size  $a$  (without replacement, treating all elements of the multiset as distinct) and let  $G = A \setminus R$ . Notice that the distribution of  $R$  and  $G$  are identical to that of  $R$  and  $G$  discussed above (namely, they both contain  $a$  independent random samples from  $X$  according to  $q$ ). We now condition on the multiset  $A = R \cup G$  and show that the event “ $E_1 \mid A$ ” happens with low probability (no matter what  $A$  is). To do so we analyze an event implied by “ $E_1 \mid A$ ” that is easier to analyze. Specifically, it holds that

$$\Pr[E_1 \mid A] \leq \Pr \left[ \exists f \in F : |\nu_{R \cup G}(Sf/s) - \nu_G(Sf/s)| \geq \frac{\varepsilon}{4} \bar{f} \mid A \right]$$

Finally, at this point we need to take the key step around the union bound in the expression above (i.e., the term “ $\exists f \in F$ ”). In the classic VC argument for binary functions this step is trivial, but in our case it is more involved and is actually what pivotally specifies our notion of an “ $\varepsilon$ -cover-code.” The remainder of the proof is devoted to this step. Details omitted.

**4.4 Proof of Theorem 4.1.** We now are ready to prove Theorem 4.1: Let  $\mathcal{X} = (\mathbb{R}^d, \varrho)$ , and  $F = W(\mathcal{X}, k, \alpha)$ . In Theorem 4.2 we show that one may assume that  $\varrho$  is  $\Gamma$ -well behaved, for  $\Gamma = (cd/\sqrt{\varepsilon})^d$  (here  $c$  is a sufficiently large constant). For a set  $A$  of size  $a$ , we have shown in Theorem 4.3 that for such well behaved  $\varrho$  the family  $F$  has cover codes of cardinality  $\left[ \left( \frac{Sa}{\varepsilon} \right)^{\Theta(\alpha^2+1)} \Gamma \right]^{2dk} = \left( \frac{cd^d Sa}{\varepsilon^d} \right)^{\Theta((\alpha^2+1)dk)}$ . As  $\left( \frac{cd^d Sa}{\varepsilon^d} \right)^{\Theta((\alpha^2+1)dk)} \leq \frac{1}{8} e^{\frac{a\varepsilon^2}{200S^2}}$  for values of  $a$  of size at least  $\Theta \left( \frac{(\alpha^2+1)d^2 k S^2}{\varepsilon^2} \log \left( \frac{(\alpha+1)dkS}{\varepsilon} \right) \right)$  by Theorem 4.4 (and a slight change of parameters to compensate for the loss of Theorem 4.2) we may conclude Theorem 4.1 stated in the beginning of this section.

## 5 Finding $\varepsilon$ -approximators for $W(\mathcal{X}, k, \alpha)$ via bi-criteria approximation

Algorithmically, our method for finding an  $\varepsilon$ -approximator described in Section 4 is a reduction to bi-criteria approximation. In the Euclidean case, such bi-criteria approximations already exist, and our results imply that further progress on such algorithms for other norms will immediately lead to efficient clustering algorithms in those norms.

Conceptually our method is simple: First, the algorithm computes a distribution  $q$  on  $X = \{x_1, \dots, x_n\}$ ; then it selects independent samples from  $X$  according to  $q$ . In what follows we show how one can construct the distribution  $q$ .

A “bi-criteria” approximation to the problem is a function  $f^* \in W(\mathcal{X}, \beta k, \alpha)$  such that  $\bar{f}^* \leq c \min_{f \in W(\mathcal{X}, k, \alpha)} \bar{f}$ . Here, both  $\beta$  and  $c$  are parameters that will affect the running time computed below.

Assuming that finding  $f^*$  can be done in time  $T_{\beta, c}$ , we show that the distribution  $q$  and an  $\varepsilon$ -approximator can be found efficiently in time  $\simeq T_{\beta, c} + O(|X|d\beta k)$ . Namely we prove Theorem 1.2 stated in the Introduction.

In general to compute  $q$  one really only needs to compute  $s(x)$  (a bound on the sensitivity) for each  $x \in X$ . Recall that in our setting it is necessary that the values  $s(x)$  we compute be greater or equal to the ideal values  $\sigma(x) = \sup_x \frac{f(x)}{\bar{f}}$  corresponding to our given family  $F$  and distribution  $\mu$ . The size of our  $\varepsilon$ -approximator resulting from  $q$  will depend on the value of  $S = \int_X s(x) d\mu$ . In what follows we show how to efficiently compute such  $s(x)$  for which  $S$  is at most  $O(2^{2\alpha}(c + \beta k))$  (no matter what  $\mu$  is). This is comparable to the (non-constructive) bounds we give on  $S$  in Section 3. Moreover, this suffices to prove Theorem 1.2. We now present our theorem for this

section.

**THEOREM 5.1.** *Let  $F = W(\mathcal{X}, k, \alpha)$ . Let  $f^* \in W(\mathcal{X}, \beta k, \alpha)$  such that  $\bar{f}^* \leq c \min_{f \in F} \bar{f}$ . Given  $f^*$ , one can compute for all  $x \in X$  a value  $s(x) \geq \sigma(x)$  and a corresponding  $q(x)$  in time  $O(|X|d\beta k)$ . The values of  $s(x)$  computed will satisfy  $S = \int_X s(x) d\mu \leq O(2^{2\alpha}(c + \beta k))$ .*

The proof of Theorem 5.1 is omitted. However, we note that its proof closely resembles that of Theorem 3.1 in which we bounded the value of  $s(x)$  (and thus  $S$ ) for families  $W(\mathcal{X}, k, \alpha)$ . More specifically, instead of setting  $f^*$  to be an optimal function in  $W(\mathcal{X}, k, \alpha)$  as in the proof of Theorem 3.1, here we set  $f^*$  to be the bi-criteria approximation.

In what follows we present the values of  $s(x)$  obtained in the proof of Theorem 5.1. Let  $f^* = f_{u_1^*, \dots, u_{\beta k}^*}$  be a function in  $W(\mathcal{X}, \beta k, \alpha)$  for which  $\bar{f}^* \leq c \min_{f \in W(\mathcal{X}, k, \alpha)} \bar{f}$ . Let  $\bar{f}^* = \Delta$ . Let  $U_i$  be the Voronoi cell of  $u_i^*$ , and let  $p_i = \mu(U_i)$ . Let  $m_i = \frac{1}{p_i} \int_{U_i} \varrho(x - u_i^*)^\alpha d\mu(x)$ , so  $\Delta = \sum p_i m_i$ . By a simple Markov inequality, for each  $i$ ,  $\mu(B(u_i^*, (2m_i)^{1/\alpha}) \cap U_i) \geq p_i/2$ . (Here  $B(x, r)$  denotes the closed ball of radius  $r$  about  $x$ .) Using the above definitions, we show

**CLAIM 1.** (a) *Let  $\alpha \leq 1$ . Setting  $s(x) = \frac{2c(2m_i + \varrho(x - u_i^*)^\alpha)}{\bar{f}^*} + \frac{4}{p_i}$  for  $x \in U_i$  satisfies  $s(x) \geq \max_f f(x)/\bar{f}$  and  $S = \int s d\mu \leq 6c + 4\beta k$ .*

(b) *Let  $\alpha \geq 1$ . Setting  $s(x) = \frac{2^{2\alpha} c m_i + 2^\alpha c \varrho(x - u_i^*)^\alpha}{\bar{f}^*} + \frac{2^{2\alpha}}{p_i}$  for  $x \in U_i$  satisfies  $s(x) \geq \max_f f(x)/\bar{f}$  and  $S = \int s d\mu \leq 2^{2\alpha}(\beta k + c) + 2^\alpha c$ .*

A small remark is in place. If the norm  $\varrho$  corresponding to  $\mathcal{X}$  is not well behaved, then naively following the flow of theorems in Section 4 one would need to find a well behaved approximation  $\varrho'$  to  $\varrho$  and use it in Theorem 5.1 above. However, this is not necessary. By multiplying the values of  $s(x)$  that correspond to  $\varrho$  computed in Theorem 5.1 by  $O((1 + \varepsilon)^\alpha)$  (to compensate for the slight difference between  $\varrho$  and  $\varrho'$  stated in Theorem 4.2) we obtain values  $s(x)$  (and thus a bound on  $S$ ) that also correspond to  $\varrho'$ .

## References

- [1] Pankaj Agarwal, Sarel Har-Peled, and Kasturi R. Varadarajan. Geometric approximation via coresets. In *Current Trends in Combinatorial and Computational Geometry*. Cambridge University Press, 2005.
- [2] N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haultner. Scale-sensitive dimensions, uniform convergence, and learnability. *J. ACM*, 44(4):615–631, July 1997.

- [3] N. Alon and B. Sudakov. On two segmentation problems. *J. of Algorithms*, 33:173–184, 1999.
- [4] P. L. Bartlett and P. M. Long. Prediction, learning, uniform convergence, and scale-sensitive dimensions. *J. Comput. System Sci.*, 56:174–190, 1998.
- [5] P. L. Bartlett, P. M. Long, and R. C. Williamson. Fat-shattering and the learnability of real-valued functions. *J. Comput. System Sci.*, 52:434–452, 1996.
- [6] S. Ben-David, N. Cesa-Bianchi, D. Haussler, and P. M. Long. Characterizations of learnability for classes of  $\{0, \dots, n\}$ -valued functions. *J. of Comput. Syst. Sci.*, 50(1):74–86, 1995.
- [7] S. Ben-David, N. Cesa-Bianchi, and P. M. Long. Characterizations of learnability for classes of  $\{0, \dots, n\}$ -valued functions. In *Proc. 5'th Annual ACM COLT*, pages 333–340, 1992.
- [8] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *J. ACM*, 36(4):929–965, October 1989.
- [9] B. Chazelle. *The Discrepancy Method*. Cambridge U. Press, 2000.
- [10] K. Chen. On coresets for  $k$ -median and  $k$ -means clustering in metric and Euclidean spaces and their applications. In *Proc. 17th SODA*, pages 1177–1185, 2006.
- [11] W. Fernandez de la Vega and C. Kenyon. A randomized approximation scheme for metric max-cut. In *39th Ann. Symp. on Foundations of Computer Science*, pages 468–471. IEEE, 1998.
- [12] R. M. Dudley. Metric entropy of some classes of sets with differentiable boundaries. *J. Approximation Theory*, 10(3):227–236, 1974.
- [13] R. M. Dudley, E. Giné, and J. Zinn. Uniform and universal Glivenko-Cantelli classes. *J. Theoretical Probability*, 4(3):485–510, 1991.
- [14] M. Effros and L. J. Schulman. Deterministic clustering with data nets. <http://eccc.hpi-web.de/eccc-reports/2004/TR04-050/index.html>, 2004. Technical Report TR04-085, Elec. Colloq. Comp. Complexity, 2004.
- [15] D. Feldman, M. Monemizadeh, and C. Sohler. A PTAS for  $k$ -means clustering based on weak coresets. In *SOCG '07: Proceedings of the twenty-third annual symposium on Computational geometry*, pages 11–18, 2007.
- [16] B. Grünbaum. Venn diagrams and independent families of sets. *Mathematics Magazine*, 48(1):12–23, 1975.
- [17] S. Har-Peled. Coresets for Discrete Integration and Clustering. In *proceedings of FSTTCS*, pages 33–44, 2006.
- [18] S. Har-Peled and A. Kushal. Smaller coresets for  $k$ -median and  $k$ -means clustering. In *SOCG '05: Proceedings of the twenty-first annual symposium on Computational geometry*, pages 126–134, New York, NY, USA, 2005. ACM.
- [19] S. Har-Peled and S. Mazumdar. Coresets for  $k$ -means and  $k$ -median clustering and their applications. In *Proc. 36th STOC*, pages 291–300, 2004.
- [20] D. Haussler and E. Welzl. Epsilon-nets and simplex range queries. *Discrete Comput. Geom.*, 2:127–151, 1987.
- [21] D. Hilbert. Über die Darstellung definiter Formen als Summe von Formenquadraten. *Math. Ann.*, 32:342–350, 1888. Ges. Abh. 2:415–436, Chelsea Publishing Co., New York, 1965.
- [22] M. J. Kearns and R. E. Schapire. Efficient distribution-free learning of probabilistic concepts. *J. Comput. System Sci.*, 48:464–497, 1994.
- [23] J. Kleinberg, C. Papadimitriou, and P. Raghavan. Approximation algorithms for segmentation problems. In *Proc. 30th STOC*, 1998.
- [24] R. Kress. *Numerical Analysis*. Springer, 1998.
- [25] M. Langberg and L. J. Schulman. Contraction and expansion of convex sets. In *proceedings of CCCG*, pages 25–28, 2004.
- [26] M. Langberg and L. J. Schulman. Universal  $\varepsilon$ -approximators for integrals. *Manuscript, available at arXiv.org*, 2009.
- [27] J. Matoušek. *Lectures on Discrete Geometry*. Springer, 2002.
- [28] D. Q. Naiman and H. P. Wynn. Independent Collections of Translates of Boxes and a Conjecture due to Grünbaum. *Discrete and Computational Geometry*, 9:101–105, 1993.
- [29] M. Naszódi and S. Taschuk. On the Transversal Number and VC-Dimension of Families of Positive Homothets of a Convex Body. *Discrete Mathematics*. To appear. Also <http://arxiv.org/abs/0907.5223>.
- [30] B. K. Natarajan. On learning sets and functions. *Machine Learning*, 4:67–97, 1989.
- [31] D. Pollard. *Convergence of stochastic processes*. Springer, 1984.
- [32] N. Sauer. On the density of families of sets. *J. Combinatorial Theory Ser. A*, 13:145–147, 1972.
- [33] L. J. Schulman. Clustering for edge-cost minimization. In *Proc. 32nd STOC*, pages 547–555, 2000.
- [34] S. Shelah. A combinatorial problem, stability and order for models and theories in infinitary languages. *Pacific J. Math.*, 41:247–261, 1972.
- [35] V. N. Vapnik. Inductive principles of the search for empirical dependences. In *Proc. 2nd Annual Workshop on Computational Learning Theory*, 1989.
- [36] V. N. Vapnik and A. Y. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–280, 1971.
- [37] Wikipedia. Sparse grid. [http://en.wikipedia.org/wiki/Sparse\\_grid](http://en.wikipedia.org/wiki/Sparse_grid). [Online; version of April 21, 2009].