

# HALFSPACE MATRICES

ALEXANDER A. SHERSTOV

**Abstract.** We introduce the notion of a *halfspace matrix*, which is a sign matrix  $A$  with rows indexed by linear threshold functions  $f$ , columns indexed by inputs  $x \in \{-1, 1\}^n$ , and the entries given by  $A_{f,x} = f(x)$ . We use halfspace matrices to solve the following problems.

In communication complexity, we exhibit a Boolean function  $f$  with discrepancy  $\Omega(1/n^4)$  under every product distribution but  $O(\sqrt{n}/2^{n/4})$  under a certain non-product distribution. This partially solves an open problem of Kushilevitz and Nisan (1997).

In learning theory, we give a short and simple proof that the statistical-query (SQ) dimension of halfspaces in  $n$  dimensions is less than  $2(n+1)^2$  under all distributions. This improves on the  $n^{O(1)}$  estimate from the fundamental paper of Blum et al. (1998). We show that estimating the SQ dimension of natural classes of Boolean functions can resolve major open problems in complexity theory, such as separating  $\text{PSPACE}^{\text{cc}}$  and  $\text{PH}^{\text{cc}}$ .

Finally, we construct a matrix  $A \in \{-1, 1\}^{N \times N^{\log N}}$  with dimension complexity  $\log N$  but margin complexity  $\Omega(N^{1/4}/\sqrt{\log N})$ . This gap is an exponential improvement over previous work. We prove several other relationships among the complexity measures of sign matrices, complementing work by Linial et al. (2006, 2007).

**Keywords.** Linear threshold functions, communication complexity, complexity measures of sign matrices, complexity of learning

**Subject classification.** 03D15, 68Q15, 68Q17.

## 1. Introduction

A *linear threshold function*, or *halfspace*, is a Boolean function  $f$  representable as  $f(x) \equiv \text{sign}(\sum_{i=1}^n a_i x_i - \theta)$  for some reals  $a_1, \dots, a_n, \theta$ . We introduce the notion of a *halfspace matrix*, which is a  $\pm 1$ -valued matrix  $A$  with rows indexed by halfspaces, columns indexed by inputs  $x \in \{-1, 1\}^n$ , and the entries given by  $A_{f,x} = f(x)$ . We demonstrate the potential of halfspace matrices in the study of complexity. Specifically, we use halfspace matrices to answer nontrivial open

questions in communication complexity, the complexity of sign matrices, and the complexity of learning.

Our work is inspired by Forster’s groundbreaking result (2002) on the sign-representation of Boolean matrices by real ones. Forster’s work has had exciting applications, including a linear lower bound on communication complexity in the unbounded-error model (Forster 2002), extremal results on Euclidean embeddings (Forster 2002; Forster *et al.* 2003; Forster & Simon 2006), and lower bounds for depth-2 threshold circuits (Forster *et al.* 2001). This paper builds on Forster’s discovery and related work to illustrate the power of halfspace matrices.

**1.1. Communication complexity.** Among the primary models in communication complexity is the *randomized model* (Kushilevitz & Nisan 1997, Chapter 3). Two parties, Alice and Bob, have access to disjoint parts  $x, y \in \{-1, 1\}^n$  of the input to a fixed function  $f : \{-1, 1\}^n \times \{-1, 1\}^n \rightarrow \{-1, 1\}$  and must therefore communicate to evaluate  $f(x, y)$ . They are allowed to use shared random bits. On every input, the players must compute the correct value with probability at least  $2/3$ . The *cost* of a protocol is number of bits exchanged in the worst case. The *randomized complexity*  $R_{1/3}^{\text{pub}}(f)$  of a function  $f$  is the cost of the best protocol for  $f$ .

The standard approach to proving lower bounds on  $R_{1/3}^{\text{pub}}(f)$  is to analyze the *distributional complexity*  $D_{1/3}^\mu(f)$  instead. One defines a probability distribution  $\mu$  on  $\{-1, 1\}^n \times \{-1, 1\}^n$  and argues that the cost  $D_{1/3}^\mu(f)$  of the best deterministic protocol with error at most  $1/3$  under  $\mu$  must be high. Yao’s well-known Minimax Principle states that  $R_{1/3}^{\text{pub}}(f) = \max_\mu \{D_{1/3}^\mu(f)\}$ . The main design question, then, is what distribution  $\mu$  to consider. While *product* distributions  $\mu(x, y) = \mu_X(x)\mu_Y(y)$  are easier to analyze, they do not always yield the optimal lower bounds. A standard example of this phenomenon is the set disjointness function DISJ: every product distribution  $\mu$  has  $D_{1/3}^\mu(\text{DISJ}) = O(\sqrt{n} \log n)$  (Kushilevitz & Nisan 1997), although  $R_{1/3}^{\text{pub}}(\text{DISJ}) = \Theta(n)$  (Kalyanasundaram & Schnitger 1992; Razborov 1992). This motivates the following intriguing question in communication complexity, due to Kushilevitz and Nisan.

**PROBLEM** (Kushilevitz & Nisan 1997, page 37). *Can restricting the distribution  $\mu$  to be a product distribution affect the resulting lower bound on  $R_{1/3}^{\text{pub}}(f)$  by more than a polynomial factor? Formally:*

$$R_{1/3}^{\text{pub}}(f) \stackrel{?}{=} \left( \max_{\mu \text{ product}} \{D_{1/3}^\mu(f)\} \right)^{O(1)}.$$

Since its formulation 10 years ago, this problem has seen little progress. The only work known to us is due to Kremer *et al.* (1999), who study the restriction of this problem to *one-way* protocols. They obtain a separation of  $O(1)$  vs.  $\Omega(n)$  for the “greater than” function GT, in the one-way model. Unfortunately, a function can have vastly different communication complexity in the one-way and normal (two-way multi-round) randomized models. Such is the case of GT, whose one-way randomized complexity is  $\Omega(n)$  but normal randomized complexity  $O(\log n)$ . Therefore, new techniques are needed to answer the Kushilevitz-Nisan question in its original formulation.

This motivates us to pursue a different approach. The chief source of lower bounds on distributional complexity  $D_{1/3}^\mu(f)$  and thus randomized complexity  $R_{1/3}^{\text{pub}}(f)$  is the so-called *discrepancy method*. The method lower-bounds  $D_{1/3}^\mu(f)$  in terms of a quantity called *discrepancy*,  $\text{disc}_\mu(f)$ . (Small discrepancy implies high communication complexity.) A natural question to ask is whether there can be a large gap between the discrepancy under product and non-product distributions. Our first main result states that, in fact, this gap can be exponential:

**THEOREM 1.1** (Discrepancy gap). *There exists an (explicitly given) function  $f : \{-1, 1\}^n \times \{-1, 1\}^{n^2} \rightarrow \{-1, 1\}$  for which  $\text{disc}_\mu(f) = \Omega(1/n^4)$  under every product distribution  $\mu$  but  $\text{disc}_\lambda(f) = O(\sqrt{n}/2^{n/4})$  under a certain non-product distribution  $\lambda$ .*

We thus establish that discrepancy-based methods can yield exponentially worse lower bounds on randomized complexity  $R_{1/3}^{\text{pub}}(f)$  if restricted to product distributions. This is the first nontrivial discrepancy gap obtained for any function.

**Recent progress.** Building on the ideas in this manuscript, the author has recently solved (Sherstov 2007) the Kushilevitz-Nisan problem in its original formulation. Namely, we have shown the existence of a function  $f : \{-1, 1\}^n \times \{-1, 1\}^{n^2} \rightarrow \{-1, 1\}$  for which  $\max_{\mu \text{ product}} \{D_\epsilon^\mu(f)\} = O(1)$  but  $R_{1/3}^{\text{pub}}(f) = \Omega(n)$ , where  $\epsilon > 0$  is an arbitrary constant. Furthermore, the function  $f$  in question has discrepancy  $O(2^{-n(1/2-\epsilon)})$ , which is almost the smallest possible. In particular,  $f$  is a hardest function for every major model of communication. Yet, the distributional method restricted to product distributions can certify at best an  $\Omega(1)$  communication lower bound for  $f$ .

The work in (Sherstov 2007) also yields an essentially optimal separation,  $\Omega(1)$  vs.  $O(2^{-n(1/2-\epsilon)})$ , between discrepancy under product and nonproduct distributions, improving on Theorem 1.1 above.

However, the results of Sherstov (2007) do not supersede Theorem 1.1. Specifically, the Kushilevitz-Nisan problem and the optimal discrepancy gap are settled there using the probabilistic method, for a *nonexplicit* function. All functions in this work, on the other hand, are given explicitly.

**1.2. Complexity of sign matrices.** A *sign matrix* is any matrix with  $\pm 1$  entries. A systematic study of sign matrices from a complexity-theoretic perspective has been recently initiated by Linial *et al.* (2006). Apart from the inherent interest of this subject as a new area of complexity, Linial *et al.* observe that several major problems in theoretical computer science are questions about sign matrices. Indeed, research into sign matrices has already yielded excellent complexity results (Forster 2002; Forster *et al.* 2001).

Our paper continues this investigation, focusing on the two main complexity measures of a sign matrix: dimension and margin complexity. Their formal definition is as follows. A *Euclidean embedding* of a sign matrix  $A \in \{-1, 1\}^{M \times N}$  is a collection of unit-length vectors  $\mathbf{u}_1, \dots, \mathbf{u}_M \in \mathbb{R}^k$  and  $\mathbf{v}_1, \dots, \mathbf{v}_N \in \mathbb{R}^k$  (for some  $k$ ) such that  $\langle \mathbf{u}_i, \mathbf{v}_j \rangle \cdot A_{ij} > 0$  for all  $i, j$ . The integer  $k$  is the *dimension* of the embedding. The quantity  $\gamma = \min_{i,j} |\langle \mathbf{u}_i, \mathbf{v}_j \rangle|$  is the *margin* of the embedding. The *dimension complexity*  $\text{dc}(A)$  is the smallest dimension of an embedding of  $A$ . The *margin complexity*  $\text{mc}(A)$  is the minimum  $1/\gamma$  over all embeddings of  $A$ .

Both dimension complexity and margin complexity have drawn much interest (Ben-David *et al.* 2003; Forster 2002; Forster *et al.* 2001, 2003; Forster & Simon 2006; Linial *et al.* 2006). In addition to their roles as complexity measures, dimension complexity is a key player in the unbounded-error model of communication (Alon *et al.* 1985; Paturi & Simon 1986), and margin complexity is the central notion in the highly successful *kernel methods* of machine learning (Burges 1998; Vapnik 1995).

Using the random projection technique of Arriaga & Vempala (2006), it is straightforward to show (Ben-David *et al.* 2003) that

$$\text{dc}(A) = O(\text{mc}(A)^2 \log(M + N))$$

for every  $M \times N$  sign matrix. This observation has had important algorithmic applications, such as the algorithm of Klivans & Servedio (2004) for learning certain intersections of halfspaces. In this paper, we ask the opposite question: *Can one place an upper bound on margin complexity in terms of dimension complexity?* A suitable upper bound of this type would establish the distribution-free weak learnability of unrestricted intersections of two halfspaces, leading to

a major breakthrough in the area (Klivans & Sherstov 2007). Unfortunately, we give a strong negative answer to this question (Theorem 1.2 below).

This problem of estimating the gap between dimension and margin complexity has been studied by several researchers. Forster *et al.* (2003) constructed a family of matrices  $A \in \{-1, 1\}^{N \times N}$  for which  $\text{dc}(A) = O(1)$  but  $\text{mc}(A) = \Theta(\log N)$ . Srebro & Shraibman (2005) amplified this separation, obtaining  $\text{dc}(A) \leq 2^p$  and  $\text{mc}(A) = (\log N)^{\Theta(p)}$  for any choice of the parameter  $1 \leq p \leq (\log N)^{1-\epsilon}$ . Our second main theorem is an exponential improvement over these results.

**THEOREM 1.2** (Margin vs. dimension). *There is an (explicitly given) matrix  $A \in \{-1, 1\}^{N \times N^{\log N}}$  for which  $\text{dc}(A) \leq \log N$  but  $\text{mc}(A) = \Omega(N^{1/4}/\sqrt{\log N})$ .*

*Note.* An exponential separation between margin and dimension complexity was independently obtained by Buhrman, Vereshchagin & de Wolf (2007) and appeared in the same conference proceedings as this paper. The proof by Buhrman *et al.* features a different matrix and completely different techniques (approximation theory and quantum communication complexity). Buhrman *et al.* phrase their result as a separation between the classes  $\text{PP}^{\text{cc}}$  and  $\text{UPP}^{\text{cc}}$  in communication complexity.

The exponential separation in Theorem 1.2 is quite close to optimal since every matrix  $A \in \{-1, 1\}^{M \times N}$  has  $1 \leq \text{dc}(A) \leq \min\{M, N\}$  and  $1 \leq \text{mc}(A) \leq \min\{\sqrt{M}, \sqrt{N}\}$  (see Section 2).

As an application of our analysis in Theorem 1.2, we consider a problem from circuit complexity (Goldmann *et al.* 1992). Fix  $d$  arbitrary functions  $f_1, \dots, f_d : \{-1, 1\}^n \rightarrow \{-1, 1\}$ . Assume that every halfspace can be computed as a majority vote of gates from among  $f_1, \dots, f_d$ . We prove that there are halfspaces that require circuits of size  $\Omega(2^{n/4}/(d\sqrt{n}))$  in this model. This generalizes the well-known fact (Håstad 1994; Siu & Bruck 1991) that some halfspaces require exponentially large weights, and complements a result due to Goldmann *et al.* (1992). See Section 6.1 for details.

We prove a number of additional results (see Section 7). In particular, we show that the standard complexity measures ( $\text{dc}(A)$ ,  $\text{mc}(A)$ , and a new complexity measure  $\text{sq}(A)$  that we introduce) form an ordered sequence that spans the continuum between  $\text{disc}^\times(A)^{-1}$  and  $\text{disc}(A)^{-1}$ . Here  $\text{disc}^\times(A)$  is the discrepancy under product distributions, and  $\text{disc}(A)$  is general discrepancy. This close interplay between linear-algebraic complexity measures ( $\text{dc}(A)$  and  $\text{mc}(A)$ ) and those from communication complexity ( $\text{disc}^\times(A)$  and  $\text{disc}(A)$ ) is further evidence that the study of sign matrices has much to contribute to complexity theory.

**1.3. Learning theory.** We adopt the *statistical query* (SQ) model of learning, due to Kearns (1993), which is a restriction of the standard PAC learning model (Valiant 1984). Fix a set  $\mathcal{C}$  of Boolean functions  $\{-1, 1\}^n \rightarrow \{-1, 1\}$  (a *concept class*) and a distribution  $\mu$  over  $\{-1, 1\}^n$ . For each choice of an unknown function  $f \in \mathcal{C}$ , the learner in the SQ model must be able to construct an approximation to  $f$  by asking queries of the form, “What is  $\mathbf{E}_{x \sim \mu} [G(x, f(x))]$ , approximately?” Here  $G : \{-1, 1\}^n \times \{-1, 1\} \rightarrow \{-1, 1\}$  is any polynomial-time computable predicate of the learner’s choosing, distinct for each query. Extensive research has established the SQ model as a powerful and elegant abstraction of learning (Blum *et al.* 1994, 2003; Kearns & Vazirani 1994; Klivans & Sherstov 2007; Yang 2005). In particular, essentially all known PAC learning algorithms can be adapted (Kearns & Vazirani 1994) to work in the SQ model. Furthermore, SQ algorithms are inherently robust to random classification noise since they use *statistics* instead of individual labeled examples.

A measure of the learning complexity of a given concept class  $\mathcal{C}$  under a given distribution  $\mu$  is its *statistical query (SQ) dimension*,  $\text{sqdim}_\mu(\mathcal{C})$ . This complexity measure captures both the running time and the sample complexity required to weakly learn  $\mathcal{C}$  in the SQ model. Informally,  $\text{sqdim}_\mu(\mathcal{C})$  is the size of the largest subset  $\mathcal{F} \subseteq \mathcal{C}$  of (almost) mutually orthogonal functions in  $\mathcal{C}$  under  $\mu$ . We put  $\text{sqdim}(\mathcal{C}) \stackrel{\text{def}}{=} \max_\mu \{\text{sqdim}_\mu(\mathcal{C})\}$ . We delay the technical details to Section 2.

Our next result concerns the concept class of halfspaces. This concept class is arguably the most studied one (Klivans *et al.* 2004; Klivans & Servedio 2004; Klivans & Sherstov 2006, 2007; Kwek & Pitt 1998; Vempala 1997) in computational learning theory, with applications in areas as diverse as data mining, artificial intelligence, and computer vision. In a fundamental paper, Blum *et al.* (1998) gave a polynomial-time algorithm for learning halfspaces in the SQ model under arbitrary distributions. It follows from their work that the SQ dimension of halfspaces is  $O(n^c)$ , where  $c > 0$  is a sufficiently large constant. We substantially sharpen this estimate:

**THEOREM 1.3** (SQ dimension of generalized halfspaces). *Fix arbitrary functions  $\phi_1, \dots, \phi_k : \{-1, 1\}^n \rightarrow \mathbb{R}$ . Let  $\mathcal{C}$  be the set of all Boolean functions  $f$  representable as  $f(x) \equiv \text{sign}(\sum_{i=1}^k a_i \phi_i(x))$  for some reals  $a_1, \dots, a_k$ . Then  $\text{sqdim}_\mu(\mathcal{C}) < 2k^2$  under all distributions  $\mu$ .*

**COROLLARY 1.4** (SQ dimension of halfspaces). *Let  $\mathcal{C}$  be the concept class of halfspaces in  $n$  dimensions. Then  $\text{sqdim}_\mu(\mathcal{C}) < 2(n + 1)^2$  under all distributions  $\mu$ .*

Prior to our work, Simon (2006, Cor. 8) proved the special case of Theorem 1.3 for  $\mu$  uniform. We generalize his result to arbitrary distributions  $\mu$ .

The SQ dimension of halfspaces is at least  $n + 1$  under the uniform distribution (consider the functions  $x_1, x_2, \dots, x_n, 1$ ). Thus, the quadratic upper bound of Corollary 1.4 is not far from optimal. In addition to strengthening the estimate of Blum *et al.*, Theorem 1.3 has a much simpler, one-page proof that builds only on Forster’s self-contained theorem (2002). The proof of Blum *et al.* relies on nontrivial notions from computational geometry and requires a lengthy analysis of robustness under noise. That said, our result gives only a nonuniform SQ algorithm for weakly learning halfspaces under an arbitrary (but fixed and known) distribution, whereas Blum *et al.* give an explicit SQ algorithm for strongly learning halfspaces under arbitrary distributions.

To best convey the importance—outside of learning theory—of studying the SQ dimension of natural classes of functions, we establish the following final result.

**THEOREM 1.5** (On the conjecture that  $\text{IP} \in \text{PSPACE}^{\text{cc}} \setminus \text{PH}^{\text{cc}}$ ). *Let  $\mathcal{C}$  be the class of functions  $\{-1, 1\}^n \rightarrow \{-1, 1\}$  computable in  $\text{AC}^0$ . If  $\text{sqdim}(\mathcal{C}) \leq O\left(2^{2^{(\log n)^\epsilon}}\right)$  for every constant  $\epsilon > 0$ , then  $\text{IP} \in \text{PSPACE}^{\text{cc}} \setminus \text{PH}^{\text{cc}}$ .*

Thus, a suitable upper bound on the SQ dimension of  $\text{AC}^0$  circuits would separate the communication-complexity analogues of the polynomial hierarchy (PH) and polynomial space (PSPACE). This latter problem is a major unresolved question in theoretical computer science that dates back to a paper by Babai *et al.* (1986). Viewed from a different standpoint, Theorem 1.5 explains the lack of progress in designing learning algorithms for  $\text{AC}^0$ : a distribution-free SQ algorithm for weakly learning  $\text{AC}^0$  in reasonable time would settle a major open question in complexity theory.

The SQ upper bound,  $\exp\left(2^{(\log n)^{o(1)}}\right)$ , for  $\text{AC}^0$  assumed in Theorem 1.5 grows faster than any quasipolynomial function but more slowly than any exponential one ( $2^{n^\epsilon}$ ). In particular, any quasipolynomial upper bound on the SQ dimension of  $\text{AC}^0$  would separate  $\text{PH}^{\text{cc}}$  and  $\text{PSPACE}^{\text{cc}}$ . At present, however, no upper bounds better than  $2^{\tilde{O}(n^{1/3})}$  are known on the SQ dimension of polynomial-size DNF formulas, let alone  $\text{AC}^0$  circuits. We hope that our observations will draw the community’s attention to the SQ dimension as an important notion in complexity theory.

**1.4. Our techniques.** A common theme of this paper is the use of halfspace matrices to answer the extremal questions at hand. Our first technical tool is

Forster’s work (2002), which we show constrains the sign patterns of halfspace matrices in an important way. These structural constraints allow us to prove an upper bound on the SQ dimension of halfspaces (Theorem 1.3). Another technical tool that we use is a theorem of Goldmann *et al.* (1992) in communication complexity. We apply it to show that halfspace matrices still do possess some structural complexity. Together, these contrasting results lead to the discrepancy gap (Theorem 1.1) and the margin-dimension gap (Theorem 1.2).

To prove Theorem 1.2, we use a technique for lower-bounding margin complexity based on communication complexity. By contrast, all previous lower bounds (Forster 2002; Forster *et al.* 2001, 2003; Linial *et al.* 2006) for explicit matrices have been based solely on linear algebra. Most of these previous techniques yield identical bounds on dimension and margin complexity, and thus cannot yield the gap of Theorem 1.2.

Finally, our proof of Theorem 1.5 regarding  $\text{PSPACE}^{\text{cc}} \setminus \text{PH}^{\text{cc}}$  builds on a manuscript of Razborov (1989) and a combinatorial observation due to Lokam (2001).

**Organization.** After some technical preliminaries, we first prove the SQ dimension upper bound (Theorem 1.3) and then use it to establish the discrepancy gap and margin-dimension gap (Theorems 1.1 and 1.2). Additional results on the complexity of sign matrices come next. We conclude the paper with observations concerning  $\text{PSPACE}^{\text{cc}} \setminus \text{PH}^{\text{cc}}$  (Theorem 1.5).

## 2. Preliminaries

This section provides relevant background on communication complexity, sign matrices, and learning theory.

**2.1. Communication complexity.** We consider Boolean functions  $f : X \times Y \rightarrow \{-1, 1\}$ . Typically  $X = Y = \{-1, 1\}^n$ , but we also allow  $X$  and  $Y$  to be arbitrary finite sets, possibly of unequal cardinality. We identify a function  $f$  with its *communication matrix*  $M = [f(x, y)]_{y, x} \in \{-1, 1\}^{|Y| \times |X|}$ . In particular, we use the terms “communication complexity of  $f$ ” and “communication complexity of  $M$ ” interchangeably (and likewise for other complexity measures, such as discrepancy). The two communication models of interest to us are the randomized model and the deterministic model, both reviewed in Section 1. The *randomized complexity*  $R_{1/2-\gamma/2}^{\text{pub}}(f)$  of  $f$  is the minimum cost of a randomized protocol for  $f$  that computes  $f(x, y)$  correctly with probability



at least  $\frac{1}{2} + \frac{\gamma}{2}$  (equivalently, with *advantage*  $\gamma$ ) for each input  $(x, y)$ . The *distributional complexity*  $D_{1/2-\gamma/2}^\mu(f)$  is the minimum cost of a deterministic protocol for  $f$  that has error at most  $\frac{1}{2} - \frac{\gamma}{2}$  (equivalently, *advantage*  $\gamma$ ) with respect to the distribution  $\mu$  over the inputs.

A distribution  $\mu$  over  $X \times Y$  is called a *product distribution* if it can be represented as  $\mu = \mu_X \times \mu_Y$  (meaning that  $\mu(x, y) = \mu_X(x)\mu_Y(y)$  for all  $x, y$ ), where  $\mu_X$  and  $\mu_Y$  are distributions over  $X$  and  $Y$ , respectively. A *rectangle* of  $X \times Y$  is any set  $R = A \times B$  with  $A \subseteq X$  and  $B \subseteq Y$ . For a fixed distribution  $\mu$  over  $X \times Y$ , the *discrepancy* of  $f$  is defined as

$$\text{disc}_\mu(f) = \max_R \left| \sum_{(x,y) \in R} \mu(x,y) f(x,y) \right|,$$

where the maximum is taken over all rectangles  $R$ . We define  $\text{disc}(f) = \min_\mu \text{disc}_\mu(f)$ . We let  $\text{disc}^\times(f)$  denote the minimum discrepancy of  $f$  under *product* distributions. Clearly,  $\text{disc}(f) \leq \text{disc}^\times(f)$ , and we will show that there can be an exponential gap between these quantities.

The *discrepancy method* is a powerful technique that lower-bounds the randomized and distributional complexity in terms of the discrepancy:

**PROPOSITION 2.1** (Kushilevitz & Nisan 1997, pp. 36–38). *For every Boolean function  $f(x, y)$ , every distribution  $\mu$ , and every  $\gamma > 0$ ,*

$$R_{1/2-\gamma/2}^{\text{pub}}(f) \geq D_{1/2-\gamma/2}^\mu(f) \geq \log_2 \frac{\gamma}{\text{disc}_\mu(f)}.$$

**2.2. Discrepancy under product distributions.** We now recall a useful technique from the literature for estimating the discrepancy under a product distribution. Our starting point is an important observation that arises as a special case in the work of Ford & Gál (2005, Thm. 3.1) on multiparty communication complexity. It is also implicit in an article by Raz (2000, Lem. 5.1).

**LEMMA 2.2** (Ford & Gál 2005; Raz 2000). *Let  $M \in \{-1, 1\}^{|X| \times |Y|}$ , and let  $\mu$  be a probability distribution over  $X \times Y$ . Then there is a choice of signs  $\alpha_x, \beta_y \in \{-1, 1\}$  for all  $x \in X, y \in Y$  such that*

$$\text{disc}_\mu(M) \leq \left| \sum_{x,y} \alpha_x \beta_y \mu(x,y) M_{xy} \right|.$$

PROOF (adapted from Raz 2000). Let  $R = A \times B$  be the rectangle for which the discrepancy is achieved. Fix  $\alpha_x = 1$  for all  $x \in A$ , and likewise  $\beta_y = 1$  for all  $y \in B$ . Choose the remaining signs  $\alpha_x, \beta_y$  independently and at random. Passing to expectations,

$$\begin{aligned} & \left| \mathbf{E} \left[ \sum_{x,y} \alpha_x \beta_y \mu(x,y) M_{xy} \right] \right| \\ &= \left| \sum_{(x,y) \in R} \underbrace{\mathbf{E}[\alpha_x \beta_y]}_{=1} \mu(x,y) M_{xy} + \sum_{(x,y) \notin R} \underbrace{\mathbf{E}[\alpha_x \beta_y]}_{=0} \mu(x,y) M_{xy} \right| \\ &= \left| \sum_{(x,y) \in R} \mu(x,y) M_{xy} \right| \\ &= \text{disc}_\mu(M). \end{aligned}$$

In particular, there exists a setting  $\alpha_x, \beta_y \in \{-1, 1\}$  for all  $x, y$  with the desired property.  $\square$

Ford and Gál used Lemma 2.2 in an elegant way to relate discrepancy to the pairwise correlations of the matrix rows:

LEMMA 2.3 (Ford & Gál 2005). *For every Boolean function  $f(x, y)$  and every product distribution  $\mu = \mu_X \times \mu_Y$ ,*

$$\text{disc}_\mu(f) \leq \sqrt{\mathbf{E}_{y,y' \sim \mu_Y} \left[ \left| \mathbf{E}_{x \sim \mu_X} [f(x,y) f(x,y')] \right| \right]}.$$

PROOF (adapted from Ford & Gál 2005). By Lemma 2.2, there is a choice of values  $\alpha_x, \beta_y \in \{-1, 1\}$  for all  $x$  and  $y$  such that

$$\text{disc}_\mu(f) \leq \left| \sum_x \sum_y \mu(x,y) \alpha_x \beta_y f(x,y) \right| = \left| \mathbf{E}_{x \sim \mu_X} \mathbf{E}_{y \sim \mu_Y} [\alpha_x \beta_y f(x,y)] \right|.$$

Thus,

$$\begin{aligned}
 \text{disc}_\mu(f)^2 &\leq \left( \mathbf{E}_{x \sim \mu_X} \mathbf{E}_{y \sim \mu_Y} [\alpha_x \beta_y f(x, y)] \right)^2 \\
 &\leq \mathbf{E}_x \left[ \left( \mathbf{E}_y [\alpha_x \beta_y f(x, y)] \right)^2 \right] \\
 &= \mathbf{E}_x \mathbf{E}_{y, y'} [\alpha_x^2 \beta_y \beta_{y'} f(x, y) f(x, y')] \\
 &\leq \mathbf{E}_{y, y'} \left[ \left| \mathbf{E}_x [f(x, y) f(x, y')] \right| \right] \quad \text{since } \alpha_x^2 = |\beta_y \beta_{y'}| = 1. \quad \square
 \end{aligned}$$

**2.3. Sign matrices.** We frequently use “generic-entry” notation to specify a matrix succinctly: we write  $A = [F(i, j)]_{i, j}$  to mean that the  $(i, j)$ th entry of  $A$  is given by the expression  $F(i, j)$ . We denote vectors by boldface letters ( $\mathbf{u}, \mathbf{v}, \mathbf{e}_i$ , etc.) and scalars by plain letters ( $u_i, x_j$ , etc.) A (*Euclidean*) *embedding* of a matrix  $A \in \{-1, 1\}^{M \times N}$  is a collection of vectors  $\mathbf{u}_1, \dots, \mathbf{u}_M \in \mathbb{R}^k$  and  $\mathbf{v}_1, \dots, \mathbf{v}_N \in \mathbb{R}^k$  (for some  $k$ ) such that  $\langle \mathbf{u}_i, \mathbf{v}_j \rangle \cdot A_{ij} > 0$  for all  $i, j$ . The integer  $k$  is the *dimension* of the embedding. The quantity

$$\gamma = \min_{i, j} \frac{|\langle \mathbf{u}_i, \mathbf{v}_j \rangle|}{\|\mathbf{u}_i\| \cdot \|\mathbf{v}_j\|}$$

is the *margin* of the embedding. The *dimension complexity*  $\text{dc}(A)$  is the smallest dimension of an embedding of  $A$ . The *margin complexity*  $\text{mc}(A)$  is the minimum  $1/\gamma$  over all embeddings of  $A$ .

Let  $\mathbf{e}_i$  denote the vector with 1 in the  $i$ th component and zeroes elsewhere. The following is a trivial embedding of a sign matrix  $A = [\mathbf{a}_1 \mid \dots \mid \mathbf{a}_N] \in \{-1, 1\}^{M \times N}$ : label the rows by vectors  $\mathbf{e}_1, \dots, \mathbf{e}_M \in \mathbb{R}^M$  and the columns by vectors  $\frac{1}{\sqrt{M}} \mathbf{a}_1, \dots, \frac{1}{\sqrt{M}} \mathbf{a}_N$ . It is easy to see that this embedding has dimension  $M$  and margin  $1/\sqrt{M}$ . By interchanging the roles of the rows and columns, we see that

$$1 \leq \text{dc}(A) \leq \min\{M, N\}, \quad 1 \leq \text{mc}(A) \leq \min\{\sqrt{M}, \sqrt{N}\}$$

for every matrix  $A \in \{-1, 1\}^{M \times N}$ . We say that a matrix  $R \in \mathbb{R}^{M \times N}$  *sign-represents* a matrix  $A \in \{-1, 1\}^{M \times N}$  if  $A_{ij} R_{ij} > 0$  for all  $i, j$ . In this case, we symbolically write  $A = \text{sign}(R)$ . Observe that the dimension complexity of a sign matrix is the minimum rank of any real matrix that sign-represents it.

The *spectral* norm of  $R \in \mathbb{R}^{M \times N}$  is defined as

$$\|R\| = \max_{\|x\|=1} \|Rx\|,$$

where the vector norm  $\|\cdot\|$  is the Euclidean norm. The *Frobenius* norm of  $R$  is defined as

$$\|R\|_F = \sqrt{\sum_{i,j} R_{ij}^2}.$$

For all  $R \in \mathbb{R}^{M \times N}$ , we have

$$\|R\|_F \geq \|R\| = \sqrt{\|RR^\top\|} = \sqrt{\|R^\top R\|}.$$

A fundamental result, due to Forster, gives a lower bound on the dimension complexity of a matrix in terms of its spectral norm:

**THEOREM 2.4** (Forster 2002). *Let  $A \in \{-1, 1\}^{M \times N}$ . Then*

$$\text{dc}(A) \geq \frac{\sqrt{MN}}{\|A\|}.$$

Using the random projection technique of Arriaga & Vempala (2006), it is straightforward to show the following relationship between dimension and margin complexity.

**PROPOSITION 2.5** (Ben-David *et al.* 2003). *Let  $A \in \{-1, 1\}^{M \times N}$ . If  $A$  has an embedding with margin  $\gamma$  (in arbitrarily high dimension), then  $A$  has an embedding with margin  $\gamma/2$  and dimension  $O(\frac{1}{\gamma^2} \log(N + M))$ . In particular,  $\text{dc}(A) \leq O(\text{mc}(A)^2 \log(N + M))$ .*

**2.4. SQ dimension.** A *concept class*  $\mathcal{C}$  is any set of Boolean functions  $\{-1, 1\}^n \rightarrow \{-1, 1\}$ . Let  $\mu$  be a probability distribution over  $\{-1, 1\}^n$ . The *statistical query (SQ) dimension* of  $\mathcal{C}$  under  $\mu$ , denoted  $\text{sqdim}_\mu(\mathcal{C})$ , is the largest  $N$  for which there are  $N$  functions  $f_1, \dots, f_N \in \mathcal{C}$  with

$$\left| \mathbf{E}_{x \sim \mu} [f_i(x) \cdot f_j(x)] \right| \leq \frac{1}{N}$$

for all  $i \neq j$ . We denote  $\text{sqdim}(\mathcal{C}) \stackrel{\text{def}}{=} \max_\mu \{\text{sqdim}_\mu(\mathcal{C})\}$ . The SQ dimension of a concept class fully characterizes its weak learnability in the statistical query model: a low SQ dimension implies an efficient weak-learning algorithm, and a high SQ dimension rules out such an algorithm; see Blum *et al.* (1994) and Yang (2005, Cor. 1). A folklore fact is that when the SQ dimension of concept class  $\mathcal{C}$  is low, it is possible to select a small number of functions that, collectively, will approximate every function in  $\mathcal{C}$ . For completeness, we state this observation with a proof.

PROPOSITION 2.6 (SQ dimension and approximation). *Let  $\text{sqdim}_\mu(\mathcal{C}) = N$ . Then there is a set  $\mathcal{H} \subseteq \mathcal{C}$  with  $|\mathcal{H}| = N$  such that each  $f \in \mathcal{C}$  has  $|\mathbf{E}_\mu[f \cdot h]| > 1/(N + 1)$  for some  $h \in \mathcal{H}$ .*

PROOF. For a set  $\mathcal{F} \subseteq \mathcal{C}$ , define

$$\gamma(\mathcal{F}) \stackrel{\text{def}}{=} \max_{f_1 \neq f_2 \in \mathcal{F}} \{|\mathbf{E}_\mu[f_1 \cdot f_2]|\},$$

the largest correlation between any two functions in  $\mathcal{F}$ . Let  $\gamma^*$  be the minimum  $\gamma(\mathcal{F})$  over all  $N$ -element subsets  $\mathcal{F} \subseteq \mathcal{C}$ . Let  $\mathcal{H}$  be a set of  $N$  functions in  $\mathcal{C}$  such that  $\gamma(\mathcal{H}) = \gamma^*$  and the number of function pairs in  $\mathcal{H}$  with correlation  $\gamma^*$  is the smallest possible (over all  $N$ -element subsets  $\mathcal{F}$  with  $\gamma(\mathcal{F}) = \gamma^*$ ).

We claim that each  $f \in \mathcal{C}$  has  $|\mathbf{E}_\mu[f \cdot h]| > 1/(N + 1)$  for some  $h \in \mathcal{H}$ . If  $f \in \mathcal{H}$ , the claim is trivially true. Thus, assume that  $f \notin \mathcal{H}$ . There are two cases to consider.

$\gamma(\mathcal{H}) \leq 1/(N + 1)$ . Then  $f$  must have correlation more than  $1/(N + 1)$  with some member of  $\mathcal{H}$ : otherwise we would have  $\gamma(\mathcal{H} \cup \{f\}) \leq 1/(N + 1)$  and  $\text{sqdim}_\mu(\mathcal{C}) \geq N + 1$ .

$\gamma(\mathcal{H}) > 1/(N + 1)$ . Again,  $f$  must have correlation more than  $1/(N + 1)$  with some member of  $\mathcal{H}$ : otherwise we could improve on the number of function pairs in  $\mathcal{H}$  with correlation  $\gamma^*$  by replacing a suitably chosen element of  $\mathcal{H}$  with  $f$ .  $\square$

When analyzing the SQ dimension of a concept class under arbitrary distributions, it is often helpful (Klivans & Sherstov 2007) to consider a modified concept class in order to keep the distribution in the analysis uniform:

PROPOSITION 2.7 (Distribution change by composition). *Let  $\mathcal{C} = \{f_1, \dots, f_t\}$  be a concept class of functions  $\{-1, 1\}^n \rightarrow \{-1, 1\}$ . Define a related class  $\mathcal{C}' = \{f_1 \circ g, \dots, f_t \circ g\}$ , where  $g : \{-1, 1\}^m \rightarrow \{-1, 1\}^n$  is an arbitrary function for some  $m$ . Then  $\text{sqdim}(\mathcal{C}) \geq \text{sqdim}_U(\mathcal{C}')$ , where  $U$  denotes the uniform distribution over  $\{-1, 1\}^m$ .*

We omit the simple proof of this fact; see Klivans & Sherstov (2007) for details.

### 3. SQ dimension of halfspaces

This section establishes an SQ upper bound for halfspaces, which plays a key role in further development.

**THEOREM 1.3** (Restated from p. 6). *Fix arbitrary functions  $\phi_1, \dots, \phi_k : \{-1, 1\}^n \rightarrow \mathbb{R}$ . Let  $\mathcal{C}$  be the set of all Boolean functions  $f$  representable as  $f(x) \equiv \text{sign}(\sum_{i=1}^k a_i \phi_i(x))$  for some reals  $a_1, \dots, a_k$ . Then  $\text{sqdim}_\mu(\mathcal{C}) < 2k^2$  under all distributions  $\mu$ .*

**COROLLARY 1.4** (Restated from p. 6). *Let  $\mathcal{C}$  be the concept class of halfspaces in  $n$  dimensions. Then  $\text{sqdim}_\mu(\mathcal{C}) < 2(n+1)^2$  under all distributions  $\mu$ .*

**PROOF** (of Theorem 1.3). We shall use the same technical tool—Forster’s work on Euclidean embeddings—as Simon (2006), who proved this claim for  $\mu$  uniform.

Let  $\mu$  be an arbitrary distribution. Assume for simplicity that  $\mu$  is rational (extension to the general case is straightforward). Then the weight  $\mu(x)$  of each point  $x$  is an integral multiple of  $1/M$ , where  $M$  is a suitably large integer.

Let  $N = \text{sqdim}_\mu(\mathcal{C})$ . Then there is a set  $\mathcal{F} \subseteq \mathcal{C}$  of  $|\mathcal{F}| = N$  functions with  $|\mathbf{E}_\mu[f \cdot g]| \leq 1/N$  for all distinct  $f, g \in \mathcal{F}$ . Consider the matrix  $A \in \{-1, 1\}^{N \times M}$  whose rows are indexed by the functions in  $\mathcal{F}$ , whose columns are indexed by inputs  $x \in \{-1, 1\}^n$  (an input  $x$  indexes exactly  $\mu(x) \cdot M$  columns), and whose entries are given by  $A = [f(x)]_{f,x}$ . By Theorem 2.4,

$$(3.1) \quad N \leq \frac{\text{dc}(A)^2 \|A\|^2}{M}.$$

We complete the proof by obtaining upper bounds on  $\text{dc}(A)$  and  $\|A\|$ .

We analyze  $\text{dc}(A)$  first. Recall that each  $f \in \mathcal{F}$  has the form  $f(x) = \text{sign}(\sum_{i=1}^k a_{f,i} \phi_i(x))$ , where  $a_{f,1}, \dots, a_{f,k}$  are reals specific to  $f$ . Therefore,

$$A = [f(x)]_{f,x} = \text{sign} \left( \left[ \sum_{i=1}^k a_{f,i} \phi_i(x) \right]_{f,x} \right) = \text{sign} \left( \sum_{i=1}^k [a_{f,i} \phi_i(x)]_{f,x} \right).$$

The last equation shows that  $A$  is sign-representable by the sum of  $k$  matrices of rank 1, i.e.,

$$(3.2) \quad \text{dc}(A) \leq k.$$

We now turn to  $\|A\|$ . The  $N \times N$  matrix  $AA^\top$  is given by

$$AA^\top = \left[ M \cdot \mathbf{E}_\mu[f \cdot g] \right]_{f,g}.$$

This means that the diagonal entries of  $AA^T$  are equal to  $M$ , whereas the off-diagonal entries do not exceed  $M/N$  in absolute value. As a consequence,

$$\begin{aligned}
 \|A\|^2 &= \|AA^T\| \\
 &\leq \|M \cdot I\| + \|AA^T - M \cdot I\| \\
 &\leq \|M \cdot I\| + \|AA^T - M \cdot I\|_F \\
 (3.3) \quad &\leq M + M\sqrt{\frac{N(N-1)}{N^2}}.
 \end{aligned}$$

Substituting the estimates (3.2) and (3.3) in (3.1) yields the inequality

$$N \leq k^2 \left( 1 + \sqrt{1 - \frac{1}{N}} \right),$$

which leads to

$$N \leq 2k^2 - \frac{1}{4}.$$

This completes the proof for rational  $\mu$ . To extend the analysis to irrational distributions  $\mu$ , one considers a rational distribution  $\tilde{\mu}$  that approximates  $\mu$  closely enough and follows the same reasoning. We omit these simple manipulations.  $\square$

**REMARK 3.4.** *An easy inspection of the proof of Theorem 1.3 reveals the following stronger result. For a distribution  $\mu$ , let  $N$  be the size of the largest set  $\{f_1, \dots, f_N\} \subseteq \mathcal{C}$  with average (not maximum!) pairwise correlation at most  $\frac{1}{N}$ , i.e.,  $\frac{1}{N(N-1)} \sum_{i \neq j} (\mathbf{E}_\mu [f_i \cdot f_j])^2 \leq \frac{1}{N^2}$ . Clearly,  $N$  is at least the SQ dimension of  $\mathcal{C}$ . Theorem 1.3 establishes an upper bound on this larger quantity:  $N < 2k^2$ .*

We will also need a version of Theorem 1.3 in a slightly different terminology.

**THEOREM 3.5** (SQ dimension and dimension complexity). *Let  $A \in \{-1, 1\}^{M \times N}$  be an arbitrary matrix. View the rows  $f_1, \dots, f_M \in \{-1, 1\}^N$  of  $A$  as Boolean functions. Then  $\text{sqdim}(\{f_1, \dots, f_M\}) < 2 \text{dc}(A)^2$ .*

In stating Theorem 3.5, we implicitly extended the notion of the SQ dimension from sets of *Boolean functions* to sets of *vectors* with  $\pm 1$  components. This extension is natural since every Boolean function can be viewed as a vector with  $\pm 1$  components, and vice versa.

#### 4. A result from communication complexity

To obtain the discrepancy gap and margin-dimension gap (Theorems 1.1 and 1.2) in the next two sections, we recall a result from communication complexity. Consider the Boolean function  $\text{GHR} : \{-1, 1\}^{4n^2} \times \{-1, 1\}^{2n} \rightarrow \{-1, 1\}$  given by

$$\text{GHR}(x, y) = \text{sign} \left( 1 + \sum_{j=0}^{2n-1} y_j \sum_{i=0}^{n-1} 2^i (x_{i,2j} + x_{i,2j+1}) \right).$$

This function was defined and studied by Goldmann *et al.* (1992) in the context of separating classes of threshold circuits. Their analysis exhibits a non-product distribution with respect to which  $\text{GHR}(x, y)$  has high distributional complexity:

**THEOREM 4.1** (Goldmann *et al.* 1992, Thm. 6 and its proof). *There is an (explicitly given) non-product distribution  $\lambda$  such that every deterministic one-way protocol for GHR with advantage  $\gamma$  with respect to  $\lambda$  has cost at least  $\log(\gamma 2^{n/2}/\sqrt{n}) - O(1)$ .*

A key consequence of Theorem 4.1 for our purposes is the following result.

**LEMMA 4.2** (Discrepancy under non-product distributions). *There is a non-product distribution  $\lambda$  for which  $\text{disc}_\lambda(\text{GHR}) = O(\sqrt{n}/2^{n/2})$ .*

**PROOF.** Consider the distribution  $\lambda$  from Theorem 4.1. Let  $R$  be a rectangle for which the discrepancy  $\text{disc}_\lambda(\text{GHR})$  is achieved:

$$\text{disc}_\lambda(\text{GHR}) = \left| \sum_{(x,y) \in R} \lambda(x, y) \text{GHR}(x, y) \right|.$$

We claim that there is a deterministic one-way protocol for  $\text{GHR}(x, y)$  with constant cost and with advantage at least  $\text{disc}_\lambda(\text{GHR})$  with respect to  $\lambda$ . Namely, define

$$a = \text{sign} \left( \sum_{(x,y) \in R} \lambda(x, y) \text{GHR}(x, y) \right), \quad b = \text{sign} \left( \sum_{(x,y) \notin R} \lambda(x, y) \text{GHR}(x, y) \right).$$



Consider the protocol  $P$  that outputs  $a$  if the input is in  $R$ , and outputs  $b$  otherwise. By definition, the advantage of  $P$  with respect to  $\lambda$  is

$$\begin{aligned} & \sum_{x,y} \lambda(x,y)P(x,y)\text{GHR}(x,y) \\ &= a \sum_{(x,y) \in R} \lambda(x,y)\text{GHR}(x,y) + b \sum_{(x,y) \notin R} \lambda(x,y)\text{GHR}(x,y) \\ &\geq \left| \sum_{(x,y) \in R} \lambda(x,y)\text{GHR}(x,y) \right| \\ &= \text{disc}_\lambda(\text{GHR}). \end{aligned}$$

But by Theorem 4.1, every one-way constant-cost protocol achieves advantage at most  $O(\sqrt{n}/2^{n/2})$ . Thus,  $\text{disc}_\lambda(\text{GHR}) = O(\sqrt{n}/2^{n/2})$ .  $\square$

## 5. Discrepancy gap

Lemma 4.2 of the previous section established that  $\text{GHR}(x,y)$  has exponentially small discrepancy under a certain *non-product* distribution. Using the SQ upper bound of Section 3, we now prove that the discrepancy of  $\text{GHR}(x,y)$  under all *product* distributions is  $\Omega(1/n^4)$ .

LEMMA 5.1 (Product distributions). *Let  $\mu = \mu_X \times \mu_Y$  be a product distribution. Then  $\text{disc}_\mu(\text{GHR}) = \Omega(1/n^4)$ .*

PROOF. For each fixed  $x$ , denote  $\text{GHR}_x(y) = \text{GHR}(x,y)$ . Since each  $\text{GHR}_x$  is a halfspace in the  $2n$  variables  $y_0, \dots, y_{2n-1}$ , Theorem 1.3 implies that

$$\text{sqdim}_{\mu_Y}(\{\text{GHR}_x\}_x) \leq \text{sqdim}_{\mu_Y}(\{\text{halfspaces in } 2n \text{ dimensions}\}) = O(n^2).$$

Thus, by Proposition 2.6, there is a set  $\mathcal{H} \subseteq \{\text{GHR}_x\}_x$  of  $|\mathcal{H}| = O(n^2)$  functions such that each  $\text{GHR}_x$  has

$$\left| \mathbf{E}_{y \sim \mu_Y} [\text{GHR}_x(y) \cdot f(y)] \right| > \frac{1}{|\mathcal{H}| + 1}$$

for some  $f \in \mathcal{H}$ .

This yields the following protocol for evaluating  $\text{GHR}(x,y)$ . Alice, who knows  $x$ , sends Bob the index of the function  $f \in \mathcal{H}$  whose correlation with  $\text{GHR}_x$  is the greatest (in absolute value). Alice additionally sends Bob the sign

$\sigma \in \{+1, -1\}$  of the correlation of  $f$  and  $\text{GHR}_x$ . This communication costs  $\lceil \log |\mathcal{H}| \rceil + 1$  bits. Bob, who knows  $y$ , announces  $\sigma \cdot f(y)$  as the output of the protocol.

For every fixed  $x$ , the described protocol achieves advantage greater than  $1/(|\mathcal{H}|+1)$  over the choice  $y$ . As a result, the protocol achieves overall advantage greater than  $1/(|\mathcal{H}|+1)$  with respect to any distribution  $\mu_X$  on the  $x$ 's. Since only  $\lceil \log |\mathcal{H}| \rceil + 2$  bits are exchanged, we obtain the sought bound on the discrepancy by Proposition 2.1:

$$\text{disc}_\mu(\text{GHR}) > \frac{1}{(|\mathcal{H}|+1) \cdot 2^{\lceil \log |\mathcal{H}| \rceil + 2}} = \Omega\left(\frac{1}{n^4}\right). \quad \square$$

Lemmas 4.2 and 4.2 immediately imply the main result of this section:

**THEOREM 1.1** (Restated from p. 3). *There exists an (explicitly given) function  $f : \{-1, 1\}^n \times \{-1, 1\}^{n^2} \rightarrow \{-1, 1\}$  for which  $\text{disc}_\mu(f) = \Omega(1/n^4)$  under every product distribution  $\mu$  but  $\text{disc}_\lambda(f) = O(\sqrt{n}/2^{n/4})$  under a certain non-product distribution  $\lambda$ .*

## 6. Margin-dimension gap

To exhibit a large gap between margin complexity and dimension complexity, we consider the function  $\text{GHR}(x, y)$  from the previous section. We first note that its dimension complexity is low.

**PROPOSITION 6.1.** *The dimension complexity of  $[\text{GHR}(x, y)]_{x, y}$  is at most  $2n + 1$ .*

**PROOF.** By definition of  $\text{GHR}$ , the sign matrix  $[\text{GHR}(x, y)]_{x, y}$  is sign-represented by the real matrix

$$M = \left[ 1 + \sum_{j=0}^{2n-1} y_j \sum_{i=0}^{n-1} 2^i (x_{i, 2j} + x_{i, 2j+1}) \right]_{x, y}.$$

It is easy to verify that  $M$  has rank at most  $2n + 1$ .  $\square$

It remains to show that the margin complexity of  $[\text{GHR}(x, y)]_{x, y}$  is high. We do so using the discrepancy estimate for  $\text{GHR}(x, y)$ . By appealing to Grothendieck's inequality and linear-programming duality, Linial & Shraibman (2006) have recently given a short, elegant proof that the margin complexity and discrepancy of a matrix are equivalent up to a small multiplicative constant:

THEOREM 6.2 (Linial & Shraibman 2006). *For every matrix  $A \in \{-1, 1\}^{M \times N}$ ,*

$$\frac{1}{8 \text{mc}(A)} \leq \text{disc}(A) \leq \frac{1}{\text{mc}(A)}.$$

Lemma 4.2 and Theorem 6.2 immediately yield an estimate of the margin complexity of  $[\text{GHR}(x, y)]_{x, y}$ .

LEMMA 6.3. *The margin complexity of  $[\text{GHR}(x, y)]_{x, y}$  is  $\Omega(2^{n/2}/\sqrt{n})$ .*

Thus, Linial and Shraibman's subtle result allows us to obtain a particularly good lower bound on the margin complexity. For completeness, we note that a slightly worse bound can be obtained using well-known and more elementary ideas relating margin complexity and discrepancy (Forster *et al.* 2001; Paturi & Simon 1986). Proposition 6.1 and Lemma 6.3 readily imply the main result of this section:

THEOREM 1.2 (Restated from p. 5). *There is an (explicitly given) matrix  $A \in \{-1, 1\}^{N \times N^{\log N}}$  for which  $\text{dc}(A) \leq \log N$  but  $\text{mc}(A) = \Omega(N^{1/4}/\sqrt{\log N})$ .*

REMARK 6.4. *While Theorem 1.2 exhibits an exponential gap between the dimension complexity and margin complexity, there can be no such gap between the dimension complexity and average squared margin. Specifically, a powerful lemma due to Forster (2002) shows that any  $k$ -dimensional embedding of a given matrix  $A \in \{-1, 1\}^{M \times N}$  can be converted into another  $k$ -dimensional embedding  $\mathbf{u}_1, \dots, \mathbf{u}_M \in \mathbb{S}^{k-1}$  and  $\mathbf{v}_1, \dots, \mathbf{v}_N \in \mathbb{S}^{k-1}$  of  $A$  that has high average squared margin:  $\frac{1}{MN} \sum_{i, j} \langle \mathbf{u}_i, \mathbf{v}_j \rangle^2 \geq \frac{1}{k}$ . (Here  $\mathbb{S}^{k-1}$  denotes the  $k$ -dimensional real unit sphere.)*

**6.1. Application: circuit complexity of halfspaces.** As an application of Lemma 6.3, we study a question from circuit complexity (Goldmann *et al.* 1992). Recall that every halfspace in  $n$  variables can be represented as

$$\text{sign}(a_1 x_1 + a_2 x_2 + \dots + a_n x_n - \theta),$$

where  $a_1, a_2, \dots, a_n, \theta$  are integers called *weights*. A fundamental fact is that there are halfspaces that require weights of magnitude  $2^{\Omega(n)}$ . This fact can be deduced by an easy counting argument, since there are  $2^{\Theta(n^2)}$  distinct halfspaces (Saks 1993). A short and simple  $2^{\Omega(n)}$  lower bound for an explicit

halfspace is due to Siu & Bruck (1991). Håstad (1994) improves on that construction, obtaining an explicit halfspace that requires weight  $2^{\Theta(n \log n)}$ . The  $2^{\Theta(n \log n)}$  lower bound is best possible for any halfspace.

Consider now a slightly modified question. Instead of expressing a halfspace as a weighted sum of the singletons  $x_1, x_2, \dots, x_n, 1$ , we get to choose an arbitrary set of Boolean functions  $f_1, \dots, f_d : \{-1, 1\}^n \rightarrow \{-1, 1\}$ . We would like to know if there is a way to choose  $d = \text{poly}(n)$  such functions such that *every* halfspace is expressible as

$$\text{sign}(a_1 f_1(x) + a_2 f_2(x) + \dots + a_d f_d(x)),$$

where  $a_1, a_2, \dots, a_d$  are integers bounded by a polynomial in  $n$ . Unfortunately, the three approaches above do not yield weight lower bounds in this more general setting. Lemma 6.3, on the other hand, yields a simple solution to the problem.

**THEOREM 6.5** (Weights of halfspaces over generalized bases). *Let  $f_1, \dots, f_d : \{-1, 1\}^n \rightarrow \{-1, 1\}$  be arbitrary functions. Assume that each halfspace in  $n$  dimensions can be expressed as  $\text{sign}(\sum_{i=1}^d a_i f_i(x))$ , where  $a_1, a_2, \dots, a_d$  are integers bounded in absolute value by  $w$ . Then  $dw \geq \Omega(2^{n/4}/\sqrt{n})$ .*

Theorem 6.5 shows a continuous trade-off between the number  $d$  of base functions and the magnitude  $w$  of the weights. In particular, the weights cannot be polynomially bounded unless there are exponentially many base functions. Goldmann *et al.* (1992, Cor. 9) proved a related result, in which the set of base functions can be arbitrarily large but they must have low randomized communication complexity (e.g., majority gates, modulo gates). Theorem 6.5 complements that result.

**PROOF** (of Theorem 6.5). It will be convenient to view  $d = d(n)$  and  $w(n) = w$  as functions of  $n$ , and set  $D = d(2n)$  and  $W = w(2n)$ . Fix the Boolean functions  $f_1, \dots, f_D$  satisfying the premise of the theorem. Consider the function

$$\text{GHR}(x, y) = \text{sign} \left( 1 + \sum_{j=0}^{2n-1} y_j \sum_{i=0}^{n-1} 2^i (x_{i,2j} + x_{i,2j+1}) \right).$$

Since for each fixed  $x$ , the function  $\text{GHR}_x(y) = \text{GHR}(x, y)$  is a halfspace in the  $2n$  variables  $y_0, \dots, y_{2n-1}$ , it is representable as a weighted sum of

$f_1(y), \dots, f_D(y)$  with integer coefficients bounded in absolute value by  $W$ . This yields the following embedding of the matrix  $[\text{GHR}(x, y)]_{x, y}$ :

$$A = \left[ \sum_{i=1}^D a_i(x) f_i(y) \right]_{x, y},$$

where each  $a_i(x)$  is an integer in  $[-W, W]$ . This embedding has margin

$$\begin{aligned} \gamma &\geq \frac{\min_{x, y} |A_{xy}|}{\max_x \left\{ \sqrt{\sum_{i=1}^D a_i(x)^2} \right\} \cdot \max_y \left\{ \sqrt{\sum_{i=1}^D f_i(y)^2} \right\}} \\ &\geq \frac{1}{\sqrt{\sum_{i=1}^D W^2} \cdot \sqrt{\sum_{i=1}^D 1^2}} \\ &= \frac{1}{DW}. \end{aligned}$$

At the same time,  $\gamma \leq O(\sqrt{n}/2^{n/2})$  by Lemma 6.3. Combining these two bounds on  $\gamma$  yields

$$DW = d(2n)w(2n) \geq \Omega\left(\frac{2^{n/2}}{\sqrt{n}}\right),$$

and thus  $dw = d(n)w(n) \geq \Omega(2^{n/4}/\sqrt{n})$ .  $\square$

## 7. Complexity measures of sign matrices: An integrated view

The results of the previous sections can be interpreted, in particular, as a study of the complexity measures of sign matrices. Our goal here is to unify them into a coherent picture and better demonstrate how they relate to previous work. We start by establishing one last fact in Section 7.1. A holistic view of this research is then given in Section 7.2.

**7.1. SQ Dimension and Discrepancy.** The SQ dimension, so far viewed as a complexity measure of *concept classes*, is just as naturally viewed as a complexity measure of *sign matrices*, as follows. Given a matrix  $A \in \{-1, 1\}^{M \times N}$ , view its rows as Boolean functions. We define the *statistical-query complexity*  $\text{sq}(A)$  of  $A$  as the SQ dimension of its rows. Formally,  $\text{sq}(A) \stackrel{\text{def}}{=} \text{sqdim}(\{f_1, \dots, f_M\})$ , where  $f_1, \dots, f_M$  are the rows of  $A$ . We prove that the SQ complexity of a matrix is essentially equivalent to the minimum discrepancy of the matrix under product distributions:

**THEOREM 7.1** (SQ complexity vs. discrepancy under product distributions).  
 Let  $A \in \{-1, 1\}^{M \times N}$ . Then

$$\sqrt{\frac{1}{2} \text{sq}(A)} < \frac{1}{\text{disc}^\times(A)} < 8 \text{sq}(A)^2.$$

Theorem 7.1 establishes a link between learning theory and communication under product distributions. Another such link was discovered earlier by Kremer *et al.* (1999), who showed that the VC dimension of a sign matrix is asymptotically equal to its maximum one-way distributional communication complexity under product distributions. The two results can be compared in a straightforward way: Theorem 7.1 relates small-bias communication to small-bias learning, whereas Kremer *et al.* relate high-accuracy communication to high-accuracy learning. Perhaps surprisingly, the proofs turn out to be quite different in the two cases.

We split the proof of Theorem 7.1 in two parts: the upper bound on  $\text{disc}^\times(A)$  and the lower bound.

**LEMMA 7.2** (Upper bound). Let  $A \in \{-1, 1\}^{M \times N}$ . Then

$$\text{disc}^\times(A) < \sqrt{\frac{2}{\text{sq}(A)}}.$$

**PROOF.** Assume  $\text{sq}(A) = d$ . Then there are  $d$  rows  $f_1, \dots, f_d \in \{-1, 1\}^N$  of  $A$  and a distribution  $\mu$  on  $\{1, \dots, N\}$  such that  $|\mathbf{E}_{x \sim \mu} [f_i(x)f_j(x)]| \leq 1/d$  for all  $i \neq j$ . Let  $U$  be the uniform distribution over the  $d$  rows  $f_1, \dots, f_d$  of  $A$ . We prove the lemma by showing that  $\text{disc}_{\mu \times U}(A) < \sqrt{2/d}$ :

$$\begin{aligned} \text{disc}_{\mu \times U}(A) &\leq \sqrt{\mathbf{E}_{i,j \sim U} [|\mathbf{E}_{x \sim \mu} [f_i(x)f_j(x)]|]} && \text{by Lemma 2.3} \\ &\leq \sqrt{\frac{1}{d} \cdot 1 + \frac{d-1}{d} \cdot \frac{1}{d}} \\ &< \sqrt{\frac{2}{d}}. \end{aligned} \quad \square$$

LEMMA 7.3 (Lower bound). *Let  $A \in \{-1, 1\}^{M \times N}$ . Then*

$$\text{disc}^\times(A) > \frac{1}{8 \text{sq}(A)^2}.$$

PROOF. The proof is closely analogous to that of Lemma 5.1; indeed, Lemma 5.1 can be deduced from this lemma. Let  $\mu \times \lambda$  be an arbitrary product distribution over  $[N] \times [M]$ . We will obtain a lower bound on  $\text{disc}_{\mu \times \lambda}(A)$  by constructing an efficient protocol for  $A$  with a suitable advantage.

Let  $\text{sq}(A) = d$ . Then by Proposition 2.6, there are  $d$  rows  $f_1, \dots, f_d \in \{-1, 1\}^N$  in  $A$  such that each of the remaining rows  $f$  has  $|\mathbf{E}_{x \sim \mu} [f(x)f_i(x)]| > 1/(d+1)$  for some  $i = 1, \dots, d$ . This yields the following protocol for evaluating  $A_{yx}$ . Bob, who knows the row index  $y$ , sends Alice the index  $i$  of the function  $f_i$  whose correlation with the  $y$ th row of  $A$  is the greatest (in absolute value). Bob additionally sends Alice the sign  $\sigma \in \{+1, -1\}$  of the correlation of  $f_i$  and the  $y$ th row of  $A$ . This communication costs  $\lceil \log d \rceil + 1$  bits. Alice, who knows the column index  $x$ , announces  $\sigma \cdot f_i(x)$  as the output of the protocol.

For every fixed  $y$ , the described protocol achieves advantage greater than  $1/(d+1)$  over the choice  $x$ . As a result, the protocol achieves overall advantage greater than  $1/(d+1)$  with respect to any distribution  $\lambda$  on the rows of  $A$ . Since only  $\lceil \log d \rceil + 2$  bits are exchanged, we obtain the sought bound on the discrepancy by Proposition 2.1:

$$\text{disc}_{\mu \times \lambda}(A) > \frac{1}{(d+1) \cdot 2^{2+\lceil \log d \rceil}} \geq \frac{1}{8d^2},$$

where the second inequality holds because  $d$  is an integer. □

Lemmas 7.2 and 7.3 immediately yield Theorem 7.1.

Since  $\text{disc}^\times(A) = \text{disc}^\times(A^\top)$ , Theorem 7.1 has the interesting corollary that the rows and columns of a matrix have the same SQ dimension, up to a polynomial factor:

COROLLARY 7.4. *Let  $A \in \{-1, 1\}^{M \times N}$ . Then*

$$\left(\frac{\text{sq}(A)}{128}\right)^{1/4} < \text{sq}(A^\top) < 128 \text{sq}(A)^4.$$

**7.2. A unified picture.** At this point, we can summarize much of this paper and relevant previous work in the following succinct diagram:

$$\frac{1}{\text{disc}^\times(A)} \stackrel{=_{\text{poly}}}{\sim} \text{sq}(A) \leq_{\text{poly}} \text{dc}(A) \leq_{\text{poly}} \text{mc}(A) \approx \frac{1}{\text{disc}(A)}$$

$\left. \begin{array}{c} \text{exponential gap} \\ \text{achievable} \end{array} \right\} \leftarrow \text{---} \rightarrow \right\}$

The purpose of this schematic is to show that the standard complexity measures ( $\text{sq}(A)$ ,  $\text{dc}(A)$ ,  $\text{mc}(A)$ ) of sign matrices form an ordered spectrum that extends from  $\text{disc}^\times(A)^{-1}$  to  $\text{disc}(A)^{-1}$ . In what follows, we let  $A \in \{-1, 1\}^{M \times N}$  be an arbitrary matrix. We shall traverse the diagram from left to right, giving precise quantitative statements.

- The smallest discrepancy of a matrix under product distributions,  $\text{disc}^\times(A)$ , and the SQ complexity of that matrix,  $\text{sq}(A)$ , are within a polynomial factor of each other:  $\Theta(\sqrt{\text{sq}(A)}) \leq \text{disc}^\times(A)^{-1} \leq \Theta(\text{sq}(A)^2)$ . We establish this fact in Theorem 7.1.
- SQ complexity puts a lower bound on dimension complexity:  $\text{dc}(A) > \sqrt{\text{sq}(A)}/2$ . We prove this relationship in Theorem 3.5.
- Dimension complexity places a lower bound the margin complexity:  $\text{mc}(A) \geq \Omega(\sqrt{\text{dc}(A)}/\log(N+M))$ . This well-known result is easy to prove using random projections (Ben-David *et al.* 2003).
- In Theorem 1.2, we show that the gap between  $\text{dc}(A)$  and  $\text{mc}(A)$  can be exponentially large. In particular, we exhibit a matrix  $A \in \{-1, 1\}^{N \times N^{\log N}}$  for which  $\text{dc}(A) \leq \log N$  but  $\text{mc}(A) = \Omega(N^{1/4}/\sqrt{\log N})$ .
- Margin complexity is within a multiplicative constant of the discrepancy:  $\text{mc}(A) \leq \text{disc}(A)^{-1} \leq 8 \text{mc}(A)$ . This is a recent result due to Linial & Shraibman (2006).

In summary, this paper refines the current understanding of the complexity measures of sign matrices by proving new relationships among them and analyzing the gaps. A particularly interesting fact is that the standard complexity measures ( $\text{sq}(A)$ ,  $\text{dc}(A)$ ,  $\text{mc}(A)$ ) form an ordered sequence that spans the continuum between product-distribution discrepancy and general discrepancy. This close interplay between linear-algebraic complexity measures ( $\text{dc}(A)$  and



$\text{mc}(A)$ ) and those from communication complexity ( $\text{disc}^\times(A)$  and  $\text{disc}(A)$ ) is further evidence that the study of sign matrices has much to contribute to complexity theory.

The only piece missing from the above schematic is the gap between  $\text{sq}(A)$  and  $\text{dc}(A)$ . This gap, stated as an open problem in the preliminary version of this work, has been recently settled by the author (Sherstov 2007). Specifically, for any constant  $\epsilon > 0$ , we have shown the existence of a matrix  $A \in \{-1, 1\}^{N \times N}$  with  $\text{sq}(A) = O(1)$  and  $\text{dc}(A) = \Omega(N^{1-\epsilon})$ . This gap is essentially the best possible by definition. It completes the above taxonomy to the following picture:

$$\frac{1}{\text{disc}^\times(A)} \stackrel{=_{\text{poly}}}{\sim} \text{sq}(A) \quad \begin{array}{c} \leftarrow \text{arbitrary gap} \\ \text{achievable} \rightarrow \end{array} \quad \begin{array}{c} \leftarrow \text{exponential gap} \\ \text{achievable} \rightarrow \end{array} \quad \text{dc}(A) \stackrel{\leq_{\text{poly}}}{\sim} \text{mc}(A) \approx \frac{1}{\text{disc}(A)}$$

## 8. Application of the SQ dimension to complexity theory

This final section demonstrates that estimating the SQ dimension of natural classes of Boolean functions is an important task in complexity theory. Specifically, we show that a suitable estimate of the SQ dimension of  $\text{AC}^0$  would solve a long-standing problem, that of separating  $\text{PH}^{\text{cc}}$  from  $\text{PSPACE}^{\text{cc}}$ .

These classes in communication complexity were introduced by Babai *et al.* (1986) by analogy with computational complexity. For our purposes, it will be convenient to view  $\text{PH}^{\text{cc}}$  and  $\text{PSPACE}^{\text{cc}}$  as classes of matrices  $M_N \in \{-1, 1\}^{N \times N}$  computed by certain circuits rather than by protocols; cf. Razborov (1989). We consider only circuits with AND, OR, NOT gates. The *inputs* to a circuit are arbitrary matrices  $A \in \{-1, 1\}^{N \times N}$  whose “−1” entries form a combinatorial rectangle; this is equivalent to requiring that  $\text{rank}(A - J) \leq 1$ , where  $J$  is the all-ones matrix. The *output* of a circuit is an  $N \times N$  sign matrix computed entry-wise from the input matrices (with the usual identification of −1 with “true,” and 1 with “false”).

**DEFINITION 8.1** (Complexity classes  $\text{PH}^{\text{cc}}$  and  $\text{PSPACE}^{\text{cc}}$ ).  *$\text{PH}^{\text{cc}}$  is the class of all matrix families  $\{M_N\}$  computable by circuits of size  $\exp((\log \log N)^{O(1)})$  and constant depth, for some choice of the input matrices.  $\text{PSPACE}^{\text{cc}}$  is the*

class of all matrix families  $\{M_N\}$  that are computable by circuits of size  $\exp((\log \log N)^{O(1)})$  and depth  $(\log \log N)^{O(1)}$ , for some choice of the input matrices.

The relevance of  $\text{PH}^{\text{cc}}$  and  $\text{PSPACE}^{\text{cc}}$  to communication complexity becomes apparent when one identifies a sign matrix with its corresponding communication problem. See Babai *et al.* (1986) for several other, equivalent definitions of these classes.

It is clear that  $\text{PH}^{\text{cc}} \subseteq \text{PSPACE}^{\text{cc}}$ , and separating these classes is a major open problem (Lokam 2001; Razborov 1989). Razborov (1989, Rem. 3) argues that “the most natural candidate for  $\text{PSPACE}^{\text{cc}} \setminus \text{PH}^{\text{cc}}$  is the INNER PRODUCT MODULO 2 predicate,” defined as

$$\text{IP}_N = [(x_1 \wedge y_1) \oplus \cdots \oplus (x_{\log N} \wedge y_{\log N})]_{x,y \in \{-1,1\}^{\log N}}.$$

Indeed, it is easy to see that  $\text{IP} \in \text{PSPACE}^{\text{cc}}$ . However, neither IP nor any other explicit family of matrices is currently known to be outside  $\text{PH}^{\text{cc}}$ . We now show that a suitable estimate of the SQ dimension of  $\text{AC}^0$  would prove the conjecture that  $\text{IP} \notin \text{PH}^{\text{cc}}$ .

**THEOREM 1.5** (Restated from p. 7). *Let  $\mathcal{C}$  be the class of functions  $\{-1, 1\}^n \rightarrow \{-1, 1\}$  computable in  $\text{AC}^0$ . If  $\text{sqdim}(\mathcal{C}) \leq O\left(2^{2^{(\log n)^\epsilon}}\right)$  for every constant  $\epsilon > 0$ , then  $\text{IP} \in \text{PSPACE}^{\text{cc}} \setminus \text{PH}^{\text{cc}}$ .*

**PROOF.** It will be convenient to prove the contrapositive: if  $\text{IP} \in \text{PH}^{\text{cc}}$ , then the SQ dimension of  $\text{AC}^0$  is at least  $2^{2^{(\log n)^\epsilon}}$  under some distribution.

We start with an insight due to Lokam (2001), restated in a terminology suitable to our proof. Let  $C$  be the assumed constant-depth circuit of size  $s = 2^{(\log \log N)^c}$  that computes  $\text{IP} \in \{-1, 1\}^{N \times N}$ . Then  $C$  has at most  $s$  inputs, which we denote by  $A_1, \dots, A_s \in \{-1, 1\}^{N \times N}$ . View the rows of the input and output matrices as Boolean functions  $\{-1, 1\}^{\log N} \rightarrow \{-1, 1\}$ . Since the “−1” entries in each  $A_1, \dots, A_s$  form a combinatorial rectangle, each  $A_i$  features at most one such row function (call it  $f_i$ ) that is not identically false. As a result, the  $r$ th row of the output matrix IP can be computed as  $C_r(f_1(x), \dots, f_s(x))$ , where  $C_r$  is a constant-depth circuit of size  $s = 2^{(\log \log N)^c}$ . See Lokam (2001) for interesting other uses of this observation.

We now return to our proof. Since the rows of IP are mutually orthogonal, the  $N$  functions  $\{C_r(f_1(x), \dots, f_s(x))\}_{r=1, \dots, N}$  that form the rows of the IP matrix are mutually orthogonal under the uniform distribution on  $x$ . By Proposition 2.7, this implies that the class of constant-depth, size- $s$  circuits

has SQ dimension at least  $N$ . Setting  $n \stackrel{\text{def}}{=} 2^{(\log \log N)^c}$ , we conclude that the class of constant-depth, size- $n$  circuits (a subclass of  $\text{AC}^0$ ) has SQ dimension at least  $2^{2^{(\log n)^{1/c}}}$ .  $\square$

## Acknowledgements

I would like to thank Anna Gál, Adam Klivans, and Sasha Razborov for helpful discussions and feedback on an earlier version of this manuscript. Thanks to Harry Buhrman, Nikolai Vereshchagin, and Ronald de Wolf for useful comments.

## References

- NOGA ALON, PETER FRANKL & VOJTECH RÖDL (1985). Geometrical Realization of Set Systems and Probabilistic Communication Complexity. In *Proc. of the 26th Symposium on Foundations of Computer Science (FOCS)*, 277–280.
- ROSA I. ARRIAGA & SANTOSH VEMPALA (2006). An algorithmic theory of learning: Robust concepts and random projection. *Mach. Learn.* **63**(2), 161–182.
- LÁSZLÓ BABAI, PETER FRANKL & JANOS SIMON (1986). Complexity classes in communication complexity theory. In *Proc. of the 27th Symposium on Foundations of Computer Science (FOCS)*, 337–347.
- SHAI BEN-DAVID, NADAV EIRON & HANS ULRICH SIMON (2003). Limitations of learning via embeddings in Euclidean half spaces. *J. Mach. Learn. Res.* **3**, 441–461.
- AVRIM BLUM, ALAN M. FRIEZE, RAVI KANNAN & SANTOSH VEMPALA (1998). A Polynomial-Time Algorithm for Learning Noisy Linear Threshold Functions. *Algorithmica* **22**(1/2), 35–52.
- AVRIM BLUM, MERRICK FURST, JEFFREY JACKSON, MICHAEL KEARNS, YISHAY MANSOUR & STEVEN RUDICH (1994). Weakly learning DNF and characterizing statistical query learning using Fourier analysis. In *Proc. of the 26th Symposium on Theory of Computing (STOC)*, 253–262.
- AVRIM BLUM, ADAM KALAI & HAL WASSERMAN (2003). Noise-tolerant learning, the parity problem, and the statistical query model. *J. ACM* **50**(4), 506–519.
- HARRY BUHRMAN, NIKOLAI K. VERESHCHAGIN & RONALD DE WOLF (2007). On Computation and Communication with Small Bias. In *Proc. of the 22nd Conf. on Computational Complexity (CCC)*, 24–32.

CHRISTOPHER J. C. BURGESS (1998). A Tutorial on Support Vector Machines for Pattern Recognition. *Data Min. Knowl. Discov.* **2**(2), 121–167.

JEFF FORD & ANNA GÁL (2005). Hadamard Tensors and Lower Bounds on Multiparty Communication Complexity. In *Proc. of the 32nd International Colloquium on Automata, Languages and Programming (ICALP)*, 1163–1175.

JÜRGEN FORSTER (2002). A linear lower bound on the unbounded error probabilistic communication complexity. *J. Comput. Syst. Sci.* **65**(4), 612–625.

JÜRGEN FORSTER, MATTHIAS KRAUSE, SATYANARAYANA V. LOKAM, RUSTAM MUBARAKZJANOV, NIELS SCHMITT & HANS-ULRICH SIMON (2001). Relations Between Communication Complexity, Linear Arrangements, and Computational Complexity. In *Proc. of the 21st Conf. on Foundations of Software Technology and Theoretical Computer Science (FST TCS)*, 171–182.

JÜRGEN FORSTER, NIELS SCHMITT, HANS ULRICH SIMON & THORSTEN SUTTORP (2003). Estimating the Optimal Margins of Embeddings in Euclidean Half Spaces. *Mach. Learn.* **51**(3), 263–281.

JÜRGEN FORSTER & HANS ULRICH SIMON (2006). On the smallest possible dimension and the largest possible margin of linear arrangements representing given concept classes. *Theor. Comput. Sci.* **350**(1), 40–48.

MIKAEL GOLDMANN, JOHAN HÅSTAD & ALEXANDER A. RAZBOROV (1992). Majority Gates vs. General Weighted Threshold Gates. *Computational Complexity* **2**, 277–300.

JOHAN HÅSTAD (1994). On the Size of Weights for Threshold Gates. *SIAM J. Discret. Math.* **7**(3), 484–492.

BALA KALYANASUNDARAM & GEORG SCHNITGER (1992). The Probabilistic Communication Complexity of Set Intersection. *SIAM J. Discrete Math.* **5**(4), 545–557.

MICHAEL KEARNS (1993). Efficient noise-tolerant learning from statistical queries. In *Proc. of the 25th Symposium on Theory of Computing (STOC)*, 392–401.

MICHAEL J. KEARNS & UMESH V. VAZIRANI (1994). *An Introduction to Computational Learning Theory*. MIT Press, Cambridge.

ADAM R. KLIVANS, RYAN O’DONNELL & ROCCO A. SERVEDIO (2004). Learning intersections and thresholds of halfspaces. *J. Comput. Syst. Sci.* **68**(4), 808–840.

ADAM R. KLIVANS & ROCCO A. SERVEDIO (2004). Learning Intersections of Halfspaces with a Margin. In *Proc. of the 17th Conf. on Learning Theory (COLT)*, 348–362.

- ADAM R. KLIVANS & ALEXANDER A. SHERSTOV (2006). Cryptographic Hardness for Learning Intersections of Halfspaces. In *Proc. of the 47th Symposium on Foundations of Computer Science (FOCS)*, 553–562.
- ADAM R. KLIVANS & ALEXANDER A. SHERSTOV (2007). Unconditional lower bounds for learning intersections of halfspaces. *Machine Learning* **69**(2–3), 97–114.
- ILAN KREMER, NOAM NISAN & DANA RON (1999). On Randomized One-Round Communication Complexity. *Computational Complexity* **8**(1), 21–49.
- EYAL KUSHILEVITZ & NOAM NISAN (1997). *Communication complexity*. Cambridge University Press, New York.
- STEPHEN KWEK & LEONARD PITT (1998). PAC Learning Intersections of Halfspaces with Membership Queries. *Algorithmica* **22**(1/2), 53–75.
- N. LINIAL, S. MENDELSON, G. SCHECHTMAN & A. SHRAIBMAN (2006). Complexity Measures of Sign Matrices. *Combinatorica* To appear. Manuscript at [http://www.cs.huji.ac.il/~nati/PAPERS/complexity\\_matrices.ps.gz](http://www.cs.huji.ac.il/~nati/PAPERS/complexity_matrices.ps.gz).
- NATI LINIAL & ADI SHRAIBMAN (2006). Learning Complexity vs. Communication Complexity. Manuscript at <http://www.cs.huji.ac.il/~nati/PAPERS/lcc.pdf>.
- SATYANARAYANA V. LOKAM (2001). Spectral Methods for Matrix Rigidity with Applications to Size-Depth Trade-offs and Communication Complexity. *J. Comput. Syst. Sci.* **63**(3), 449–473.
- RAMAMOCHAN PATURI & JANOS SIMON (1986). Probabilistic communication complexity. *J. Comput. Syst. Sci.* **33**(1), 106–123.
- RAN RAZ (2000). The BNS-Chung criterion for multi-party communication complexity. *Comput. Complex.* **9**(2), 113–122.
- ALEXANDER A. RAZBOROV (1989). Ob ustoichivyh matritsah. Research report, Steklov Mathematical Institute, Moscow, Russia. In Russian. *Engl. title*: “On rigid matrices”.
- ALEXANDER A. RAZBOROV (1992). On the distributional complexity of disjointness. *Theor. Comput. Sci.* **106**(2), 385–390.
- MICHAEL E. SAKS (1993). Slicing the hypercube. *Surveys in combinatorics, 1993* 211–255.
- ALEXANDER A. SHERSTOV (2007). Communication Complexity under Product and Nonproduct Distributions. In *Electronic Colloquium on Computational Complexity (ECCC)*. Report TR07-072.

HANS ULRICH SIMON (2006). Spectral Norm in Learning Theory: Some Selected Topics. In *Proc. of the 17th Conf. on Algorithmic Learning Theory (ALT)*, 13–27.

KAI-YEUNG SIU & JEHOASHUA BRUCK (1991). On the Power of Threshold Circuits with Small Weights. *SIAM J. Discrete Math.* **4**(3), 423–435.

NATHAN SREBRO & ADI SHRAIBMAN (2005). Rank, Trace-Norm and Max-Norm. In *Proc. of the 18th Conf. on Learning Theory (COLT)*, 545–560.

LESLIE G. VALIANT (1984). A theory of the learnable. *Commun. ACM* **27**(11), 1134–1142.

VLADIMIR N. VAPNIK (1995). *The nature of statistical learning theory*. Springer-Verlag, New York.

SANTOSH VEMPALA (1997). A Random Sampling based Algorithm for Learning the Intersection of Halfspaces. In *Proc. of the 38th Symposium on Foundations of Computer Science (FOCS)*, 508–513.

KE YANG (2005). New lower bounds for statistical query learning. *J. Comput. Syst. Sci.* **70**(4), 485–509.

Manuscript received 28 August 2007

ALEXANDER A. SHERSTOV  
Department of Computer Sciences  
The University of Texas at Austin  
1 University Station C0500  
Austin, TX 78712-0233 USA  
[sherstov@cs.utexas.edu](mailto:sherstov@cs.utexas.edu)  
<http://www.cs.utexas.edu/~sherstov>