# Tight Hardness Results for Minimizing Discrepancy

Moses Charikar[*]  Alantha Newman[†]  Aleksandar Nikolov[‡]

January 13, 2011

## Abstract

In the DISCREPANCY problem, we are given $M$ sets $\{S_1, \ldots, S_M\}$ on $N$ elements. Our goal is to find an assignment $\chi$ of $\{-1, +1\}$ values to elements, so as to minimize the maximum discrepancy $\max_j |\sum_{i \in S_j} \chi(i)|$. Recently, Bansal gave an efficient algorithm for achieving $O(\sqrt{N})$ discrepancy for any set system where $M = O(N)$ [Ban10], giving a constructive version of Spencer's proof that the discrepancy of any set system is at most $O(\sqrt{N})$ for this range of $M$ [Spe85].

We show that from the perspective of computational efficiency, these results are tight for general set systems where $M = O(N)$. Specifically, we show that it is NP-hard to distinguish between such set systems with discrepancy zero and those with discrepancy $\Omega(\sqrt{N})$. This means that even if the optimal solution has discrepancy zero, we cannot hope to efficiently find a coloring with discrepancy $o(\sqrt{N})$. We also consider the hardness of the DISCREPANCY problem on sets with bounded shatter function, and show that the upper bounds due to Matoušek [Mat95] are tight for these sets systems as well.

The hardness results in both settings are obtained from a common framework: we compose a family of high discrepancy set systems with set systems for which it is NP-hard to distinguish instances with discrepancy zero from instances in which a large number of the sets (i.e. constant fraction of the sets) have non-zero discrepancy. Our composition amplifies this zero versus non-zero gap.

## 1 Introduction

In the DISCREPANCY problem, we are given $M$ sets $\{S_1, \ldots, S_M\}$ on $N$ elements. Our goal is to find an assignment $\chi$ of $\{-1, +1\}$ values (sometimes referred to as a coloring) to elements, so as to minimize the maximum discrepancy $\max_j |\sum_{i \in S_j} \chi(i)|$. Questions about the discrepancy of various set systems arise in several different areas of mathematics and theoretical computer science and have given rise to a rich body of research on the subject. For a comprehensive introduction to discrepancy and its applications, the reader is referred to [BS96, CMS95, Mat10].

A celebrated result of Spencer [Spe85] shows that any system of $M = O(N)$ sets has discrepancy at most $O(\sqrt{N})$. Spencer's proof is non-constructive and until very recently, an efficient algorithm to construct a low discrepancy coloring was not known. In a recent breakthrough, Bansal [Ban10] gave an efficient algorithm to find a low discrepancy coloring. For a given set system of $M = O(N)$, the discrepancy of the coloring produced is at most $O(\sqrt{N})$ giving a constructive version of Spencer's proof.

Thus Bansal's algorithm gives a coloring with discrepancy that matches (within constants) the worst case discrepancy for all set systems with $M = O(N)$ sets. This leaves open the tantalizing question: Can we achieve discrepancy bounds tailored to the optimal discrepancy of the input instance instead of the worst case discrepancy over all instances? In particular, can we get better guarantees for discrepancy if the optimal discrepancy for the input instance is small? Given that the existence of an efficient algorithm for achieving worst case discrepancy was open until recently, it is not surprising that very little is known about these questions. In fact Bansal mentions in his paper that even if the optimal discrepancy is *zero*, we do not know if we can efficiently find a coloring with discrepancy better than the worst case discrepancy.

In this paper, we show strong hardness results that rule out any better discrepancy guarantees for efficient algorithms. We show that from the perspective of computational efficiency, Bansal's results are tight for general set systems where $M = O(N)$. Specifically, our main theorem states:

**Theorem 1 (Main Theorem)** *It is NP-hard to distinguish between set systems on $N$ elements with discrepancy zero and those with discrepancy $\Omega(\sqrt{N})$.*

This means that even if the optimal solution has discrepancy zero, we cannot hope to efficiently find a coloring with discrepancy $o(\sqrt{N})$. Our approach to proving the Main Theorem is to compose a family of high discrepancy set systems with a family for which it is NP-hard to distinguish instances with discrepancy zero from instances in which a constant fraction of the sets have discrepancy $\Omega(1)$. Our composition amplifies this zero versus $\Omega(1)$ gap.

Our methods are general enough that we also obtain a similar theorem for set systems with bounded shatter function. For such set systems, we show that the upper bounds due to Matoušek [Mat95] are tight. The proof for this latter result involves using high discrepancy set systems that have bounded shatter function in the composition, and proving that the resulting set system also has bounded shatter function. Thus, our methods suggest a general framework where we can obtain computational lower bounds for computing the discrepancy on other restricted set systems. In particular, our composition consists of two main steps that need to be tailored to a specified type of restricted set system: (i) plug in a lower bound (i.e. high discrepancy) instance for a set system with certain specified properties, and (ii) show that the final set system maintains these specified properties. If these two steps can be carried out, the discrepancy of the lower bound instance will be translated to computational hardness of distinguishing between discrepancy zero and discrepancy equal to that of the lower bound instance.

# 2 Preliminaries

We follow the notation from Chazelle [Cha91] and use the $M \times N$ 0-1 matrix $A$ to represent a specified set system on $N$ elements; each row of the matrix $A$ is the indicator vector for a set. The discrepancy of a set system $A$, also denoted by $D_\infty(A)$, can be expressed as:

$$D_\infty(A) \quad = \quad \min_{x \in \{-1,1\}^N} ||Ax||_\infty. \tag{1}$$

Further following Chazelle, we have:

$$D_2^2(A) \quad = \quad \min_{x \in \{-1,1\}^N} ||Ax||_2^2. \tag{2}$$

The following fact exhibits a useful relationship between (1) and (2):

**Fact 1** $D_\infty(A) \geq \sqrt{\frac{D_2^2(A)}{M}}$.

The proof of our Main Theorem (Theorem 1) follows from Fact 1 together with the following stronger theorem:

**Theorem 2** *Given an $M \times N$ 0-1 matrix $A$ with $M = O(N)$, it is NP-hard to distinguish between the cases* **(1)** $D_2^2(A) = 0$, *and* **(2)** $D_2^2(A) \geq \Omega(N^2)$.

## 2.1 Overview

Our proof of Theorem 2 uses a reduction from the MAX-2-2-SET-SPLITTING problem, which is defined in Section 2.2. The reduction uses two key steps to construct a set system: (i) use of a Hadamard matrix (i.e. a lower bound instance) to design a new set system containing multisets, and (ii) preprocessing the instance of MAX-2-2-SET-SPLITTING so that the procedure in step (i) results in a set system without multisets. Hadamard set systems on $N$ elements and $N$ sets are a lower bound instance, since they have discrepancy $\Omega(\sqrt{N})$ [Cha91]. In Section 3, we illustrate step (i), by proving a variant of Theorem 2 on multisets. In Section 4, we prove Theorem 2 by showing how to create a set system in which each element appears at most once per set. Our methods are not tailored specifically to Hadamard set systems. Thus, in Section 5, we show that we can get similar hardness results for set systems with bounded shatter function by substituting lower bound instances for these set systems in place of Hadamard set systems.

## 2.2 Max-2-2-Set-Splitting

Our proofs use a reduction from the MAX-2-2-SET-SPLITTING problem, which we now define. In an instance of the MAX-2-2-SET-SPLITTING problem, we are given $m$ sets, $\mathcal{C} = \{C_1, C_2, \ldots C_m\}$, on $n$ elements, each set consisting of exactly four distinct elements. The objective is to assign each of the $n$ elements a value from $\{-1, 1\}$ so as to maximize the number of sets in which the values of the elements sum to 0. Such a set is called *2-2-split* or *split*, for short. A set whose value does not sum to 0 is called *unsplit*. Unsplit sets have possible values $\{\pm 2, \pm 4\}$. It is NP-hard to

distinguish between instances of Max-2-2-Set-Splitting for which there is an assignment such that all sets are split and instances for which any assignment will result in at least $\phi m$ unsplit sets, for some positive constant $\phi$ [Gur03]. We refer to the former and latter instances as *splittable* and *unsplittable*, respectively. (Note that this immediately implies that it is NP-hard to distinguish between set systems with discrepancy zero and those with discrepancy two, which, to the best of our knowledge, was the strongest hardness result known for the Discrepancy problem until now.) Furthermore, as noted in [CGW05], we can assume that each of the $n$ elements appears in no more than $b$ sets for some constant $b$. Thus, $m \leq bn/4 = \Theta(n)$. In order to distinguish between the sets in a given instance of Max-2-2-Set-Splitting and the sets of elements in the set systems that we construct, we refer to the former as "clauses" from hereon in.

Finally, when necessary, we can assume that $m$ is a power of two. Note that if $m$ is not a power of two, we can add at most $m$ clauses on at most $4n$ elements so that the number of clauses is now a power of two. Each of the new clauses will contain its own four distinct elements and thus these new clauses will be always be splittable. (Actually, each new clause in the "padding" need only contain two elements, since this is sufficient for the clause to have value zero.) Thus, we obtain a modified instance of Max-2-2-Set-Splitting with $n' = O(n)$ elements and $m' = O(n)$ clauses such that $m'$ is a power of two, such that the instance is unsplittable iff the original instance was unsplittable and has $\gamma m'$ unsplit clauses for any assignment of the elements for some constant $\gamma$ if the original instance was unsplittable.

## 3   Hadamard Set Systems

In this section, we illustrate one of the two key steps in our proof of Theorem 2: the use of Hadamard matrices to construct a set system. We describe this technique in the proof of the following theorem, which is an analogue of Theorem 2 for multisets.

**Theorem 3** *Given an $M \times N$ matrix $B$ with $M = O(N)$ and entries in $\mathbb{N}^+ \cap [0, b]$, where $b$ is a constant, it is NP-hard to distinguish between the cases* **(1)** $\exists y \in \{-1, 1\}^N$ *for which* $||By||_2^2 = 0$, *and* **(2)** $\forall y \in \{-1, 1\}^N$, $||By||_2^2 \geq \Omega(N^2)$.

**Proof:**   We begin with an instance of Max-2-2-Set-Splitting consisting of $n$ elements and $m$ clauses, $\{C_1, C_2, \ldots C_m\}$, of four distinct elements each. Let us represent the $n$ elements as $\mathcal{E} = \{e_1, e_2, \ldots e_n\}$. We construct a new set system, $B$, of $m$ multisets on the same $n$ elements, in which each new multiset contains $\Omega(n)$ elements. Let $H$ be an $m \times m$ Hadamard matrix and let $J$ be an $m \times m$ matrix of all ones. Let $W = \frac{1}{2}(H + J)$. To create the $i^{th}$ (multi)set in $B$, we consider the $i^{th}$ row of $W$ and include all elements from the clause $C_j$ such that $W_{ij} = 1$ (i.e. $j^{th}$ entry of the $i^{th}$ row of $W$ is a '1').

Suppose we fix an assignment, $\chi$, of the $n$ elements in $\mathcal{E}$ to values in $\{-1, 1\}$. Let $x$ be a vector of length $m$ such that $x_j$ represents the sum of the values in clause $C_j$, i.e. $x_j = \sum_{e_i \in C_j} \chi(e_i)$ and $x_j \in \{0, \pm 2, \pm 4\}$. Let $B$ (an $m \times n$ matrix) denote the system of $m$ multisets on $n$ elements such that the $i^{th}$ entry of row $j$ denotes the multiplicity of element $e_i$ in the set represented by vector $B_j$–the $j^{th}$ row of matrix $B$. Let $y$ denote the vector of length $n$ corresponding to the assignment

4

$\chi$ for the $n$ elements. Note that:

$$||By||_2^2 \quad = \quad ||Wx||_2^2. \tag{3}$$

If the assignment $\chi$ splits the original instance of MAX-2-2-SET-SPLITTING, then $||x||_2^2 = 0$ and therefore $||Wx||_2^2 = 0$. However, if $\chi$ is an unsplittable assignment, resulting in $\phi m$ of the original sets being unsplit, then $||Wx||_2^2 = \Omega(m^2) = \Omega(n^2)$. This latter fact follows from pages 18–19 of [Cha91], which we reproduce here for the sake of completeness. In this proof, we use the property that $\sum_{i>1} x_i^2 = \Omega(m)$. $\qquad\square$

**Lemma 1** *Let $W = \frac{1}{2}(H + J)$, where $H$ is a $k \times k$ Hadamard matrix and $J$ is a $k \times k$ matrix of all ones. Let $x \in \mathbb{R}^k$ be a vector such that $\sum_{i>1} x_i^2 = \Omega(k)$. Then $||Wx||_2^2 = \Omega(k^2)$.*

**Proof:** The proof of this lemma follows directly from the arguments on pages 18–19 of [Cha91] with very minor modifications. Following [Cha91], we have:

$$||Wx||_2^2 \quad = \quad (Wx)^T(Wx) \quad = \quad \frac{1}{4} \sum_{i,j} x_i x_j (H_i + J_i)^T (H_j + J_j).$$

Breaking this down, and further following [Cha91], we have:

1. $\sum_{i,j} x_i x_j H_i^T H_j = \sum_i x_i^2 H_i^T H_i$.

2. Since $J_j = H_1$, we have $\sum_{i,j} x_i x_j H_i^T J_j = \left( \sum_j x_j \right) \sum_i x_i H_i^T H_1$.

   By orthogonality, this is $(\sum_j x_j) x_1 k$. We find the same value for $\sum_{i,j} x_i x_j J_i^T H_j$.

3. The term $\sum_{i,j} x_i x_j J_i^T J_j = (\sum_i x_i)^2 k$.

Combining these terms, we have:

$$
\begin{aligned}
4||Wx||_2^2 \quad &= \quad k \sum_i x_i^2 + 2k \left( \sum_i x_i \right) x_1 + k \left( \sum_i x_i \right)^2 \\
&= \quad k \left( x_1 + \sum_i x_i \right)^2 + k \sum_{i>1} x_i^2 \\
&\geq \quad k \sum_{i>1} x_i^2 \\
&= \quad \Omega(k^2).
\end{aligned}
$$

This concludes the lemma. $\qquad\square$

# 4 Proof of Theorem 2

We now prove Theorem 2, which is essentially Theorem 3 except that the sets in the set system are no longer allowed to be multisets. This presents a difficulty that we circumvent by dividing the clauses of the original instance of Max-2-2-Set-Splitting into sets of clauses such that no pair of clauses contain the same element. In other words, each element appears in at most one clause per set. Then we run the construction from Section 3 to create a Hadamard set system on each set of clauses.

As before, we begin with an instance $\mathcal{C}$ of Max-2-2-Set-Splitting consisting of $n$ elements and $m$ clauses. Instead of creating a Hadamard set system directly from these $m$ clauses, we first divide the $m$ clauses into a constant number of sets such that in any set, each element appears in at most one clause.

**Lemma 2** *The clauses from an instance $\mathcal{C}$ of* Max-2-2-Set-Splitting *can be divided into $c$ sets, where $c$ is a constant, such that in each set, each element appears in at most one clause.*

**Proof:** Construct a graph $G = (V, E)$, where each vertex in $V$ is associated with one of the $m$ clauses in $\mathcal{C}$ and there is an edge between two vertices in $V$ if the associated pair of clauses has at least one element in common. Each element in $\mathcal{C}$ appears in at most $b$ clauses, so $|V| = m \leq bn/4$. Each vertex in $V$ has degree at most $4b$. Thus, we can color this graph with $4b + 1$ colors to obtain $c \leq 4b + 1$ "independent" sets of clauses. $\qquad\square$

**Lemma 3** *If $\mathcal{C}$ is an* unsplittable *instance of* Max-2-2-Set-Splitting*, then for each of the $2^n$ assignments to the $n$ elements, at least one of the $c$ independent sets of clauses contains at least $\Omega(n)$ unsplit clauses.*

**Proof:** Any assignment to the $n$ elements results in $\gamma m$ unsplit clauses for some constant $\gamma > 0$. At least one of the independent sets of clauses contains at least $\gamma m/c = \Omega(n)$ unsplit clauses. $\qquad\square$

Now we use Lemma 2 to construct a new set system in the following way. Beginning with an instance of Max-2-2-Set-Splitting, we can divide the clauses into $c$ sets such that the clauses in each set are "independent". Call these sets $\{T_1, T_2, \ldots T_c\}$. Let us fix an integer $h \in [1, c]$ and consider the independent set $T_h$ containing $k'$ clauses. We construct a set system on $k$ sets, where $k$ is the nearest power of two such that $k \geq k'$. (We use the padding "trick" discussed in Section 2.2.)

To do this, we use the incidence matrix $W^{(h)} = \frac{1}{2}(H + J)$, where $H$ is a $k \times k$ Hadamard matrix and $J$ is a $k \times k$ matrix of all ones. Consider the set of clauses $T_h$. For each row of $W^{(h)}$, we include all elements of the $j^{th}$ clause in $i^{th}$ set if $W_{ij}^{(h)}$ is a '1'. Note that each new set corresponding to a row of the matrix $W^{(h)}$ is not a multiset, since all elements in $T_h$ appear exactly once. Let the final set system be the union of the $c$ set systems obtained via the above procedure. The resulting set system is on $N = \Theta(n)$ elements and contains $M = O(cn)$ sets.

Now let the matrix $A$ denote this set system of $M$ sets on $N$ elements. In other words, $A$ is a $M \times N$ 0-1 matrix. Furthermore, let us consider an assignment, $\chi$, to the $N$ distinct elements. Consider the set $T_h = \{C_1, C_2, \ldots C_k\}$. Let the vector $x^{(h)} \in \{\pm 4, \pm 2, 0\}^k$ be defined as follows.

For each $C_j \in T_h$, we have:

$$x_j^{(h)} = \sum_{e_i \in C_j} \chi(e_i). \tag{4}$$

Let $y \in \{-1, 1\}^N$ denote the assignment $\chi$ for $N$ elements. Note that:

$$\|Ay\|_2^2 = \sum_{h=1}^{c} \|W^{(h)} x^{(h)}\|_2^2. \tag{5}$$

If the original instance of MAX-2-2-SET-SPLITTING is splittable, then there is an assignment $\chi$ on the $N$ elements such that $\|Ay\|_2^2 = 0$. However, if the original instance of MAX-2-2-SET-SPLITTING is unsplittable, then by Lemma 3, for at least one of the vectors, say $x^{(h)}$, it is the case that $\|x^{(h)}\|_2^2 = \Omega(N)$, and in fact $\sum_{j>1}(x_j^{(h)})^2 = \Omega(N)$. Thus, by Lemma 1, for this case, we have:

$$\sum_{h=1}^{c} \|W^{(h)} x^{(h)}\|_2^2 = \Omega(N^2). \tag{6}$$

This completes the proof of Theorem 2.

Finally, we note that the following corollary holds.

**Theorem 4** *Given an $M \times N$ 0-1 matrix $A$ for $M = O(N)$ such that each column of $A$ has at most $t$ 1's, it is NP-hard to distinguish between the cases* **(1)** $D_\infty(A) = 0$, *and* **(2)** $D_\infty(A) \geq \Omega(\sqrt{t})$.

**Proof:** If $t = O(N)$, then this theorem follows from Theorem 2 and from the fact that in the set systems constructed in the proof of Theorem 2, there are $N$ elements and $M = O(N)$ sets. Thus, each element has degree $t = O(N)$, since it can appear at most once in each set.

However, note that we can actually prove this theorem for any value of $t$ such that $\Omega(1) \leq t < \Omega(N)$. This is because after we take the independent sets of clauses from the set-splitting instance (using Lemma 2), instead of obtaining $4b + 1$ (where $b$ is a constant) independent sets, we can obtain $k$ independent sets, each with at most $t$ elements (say, between $t$ and $2t$). The number of such sets is $k = \Omega(N)/t$. We use these $k$ sets to create $k$ Hadamard set systems as described in Section 3. The maximum degree of a set over all these set systems is $\Omega(t)$. The number of unsplit clauses is $\Omega(N)$. Thus, one of the $k$ sets contains $\Omega(N)/k = \Omega(t)$ unsplit sets, and therefore has discrepancy $\Omega(\sqrt{t})$. □

## 5  Set Systems with Bounded Shatter Function

For some special classes of set systems there exist bounds that improve on the guarantees of Spencer's theorem. For example, Matoušek [Mat95] showed improved discrepancy bounds for set systems whose shatter function is polynomially bounded. Such set systems arise frequently in computational geometry and computational learning theory. Moreover, Matoušek's bounds can be made constructive using Bansal's method [Ban10]. In this section, we show tight inapproximability results for the discrepancy of set systems with polynomially bounded shatter function. They are proved using the same approach that was used for proving Theorem 2.

## 5.1  Preliminaries

Let $(X, \mathcal{S})$ be a set system on $N = |X|$ elements and $M = |\mathcal{S}|$ sets. Given $Y \subseteq X$, the trace of $Y$ on $\mathcal{S}$ is defined as $\mathcal{S}_{|Y} = \{S \cap Y : S \in \mathcal{S}\}$. The *primal shatter function* is

$$\pi_{\mathcal{S}}(s) = \max_{Y : |Y| = s} |\mathcal{S}_{|Y}|. \tag{7}$$

Matoušek proved that for set systems $(X, \mathcal{S})$ such that $\pi_{\mathcal{S}}(s) = O(s^d)$, $D_{\infty}(\mathcal{S}) = O(N^{1/2 - 1/2d})$ [Mat95]. The proof relies on the entropy lemma; since Bansal gives a constructive version of the lemma [Ban10], Maroušek's bound can be proved constructively as well. We show that this is essentially best possible.

**Theorem 5** *Given a set system $(X, \mathcal{S})$, with $|X| = N$ and $\pi_{\mathcal{S}}(s) = O(s^d)$, it is NP-hard to distinguish between the cases **(1)** $D_{\infty}(\mathcal{S}) = 0$, and **(2)** $D_{\infty}(\mathcal{S}) = \Omega(N^{1/2 - 1/2d})$.*

Recall that one of the main ingredients in the proof of Theorem 2 was a family of high discrepancy set systems: the Hadamard set systems, which are a tight example for Spencer's theorem. Analogously, in the proof of Theorem 5 we use a family of high discrepancy set systems with polynomially bounded shatter function. The family consists of systems of sets defined by halfspaces. The discrepancy lower bound for such set systems was proved by Alexander [Ale90]. We present the result as it appears in Chazelle [Cha91]. We need an extension of the original result that follows from the proof technique introduced in [CMS95].

We first need to introduce a new definition. For a set $X$ of points in $\mathbb{R}^d$, let $d_{max}(X) = \max_{x,y \in X} ||x - y||_2^2$, and, similarly, $d_{min}(X) = \min_{x,y \in X} ||x - y||_2^2$.

**Definition 1** *A set $X$ of $N$ points in $\mathbb{R}^d$ is* well-spread *if $d_{max}(X)/d_{min}(X) = O(N^{1/d})$.*

Observe that the set of vertices of a grid inside a $d$-dimensional cube is well-spread. We also note the following simple fact that will be useful in the proof of Theorem 5.

**Lemma 4** *Let $X$ be a well-spread set of $N$ points in $\mathbb{R}^d$. If $Y \subseteq X$ and $|Y| = \Omega(N)$, then $Y$ is well-spread.*

**Proof:** Since $d_{max}(Y) \leq d_{max}(X)$ and $d_{min}(Y) \geq d_{min}(X)$, $d_{max}(Y)/d_{min}(Y) \leq d_{max}(X)/d_{min}(X) = O(N^{1/d})$. By $|Y| = \Omega(N)$ we have $O(N^{1/d}) = O(|Y|^{1/d})$ and this completes the proof. $\square$

We can now state the generalized version of Alexander's lower bound.

**Lemma 5** *Let $X$ be a well-spread set of $N$ points in $\mathbb{R}^d$. For any assignment $\chi : X \to \{\pm 1, \pm 2\}$, there exists a hyperplane $H$, such that*

$$\sum_{x \in X \cap H^+} \chi(x) = \Omega(N^{1/2 - 1/2d}), \tag{8}$$

*where $H^+$ is the closed halfspace above $H$.*

8

**Proof:** If $\chi : X \to \{\pm 1\}$, then this theorem is Theorem 3.9 in [Cha91]. However, since we need $\chi : X \to \{\pm 1, \pm 2\}$, we need to modify the proof slightly. See Appendix A for details. $\qquad\square$

It is a well known fact that a set system $(X, \mathcal{S})$ of halfspaces in $\mathbb{R}^d$ has $\pi_{\mathcal{S}}(s) = O(s^d)$ (see e.g. [Mat10]). Thus, such set systems are a tight example for Matoušek's upper bound.

Our proof of the hardness of approximating discrepancy on set systems with polynomially bounded shatter function follows the structure of the proof of Theorem 2. The two key steps in the proof of Theorem 5 are using systems of halfspaces instead of Hadamard set systems, and showing that the shatter function of the final construction is bounded by $O(s^d)$.

## 5.2 Proof of Theorem 5

Once again we begin with an instance of MAX-2-2-SET-SPLITTING with $n$ elements and $m$ clauses. As in the proof of Theorem 2, we partition the clauses of the set splitting instance into sets $T_1, \ldots, T_c$, such that for any $h$ no two clauses in $T_h$ share an element. Also, if each element in the set splitting instance appears in at most $b$ clauses, then $c \leq 4b + 1$.

We construct a set system $\mathcal{S}_h$ for each $T_h$, and the final systems $\mathcal{S}$ will be the union of $\mathcal{S}_1, \ldots, \mathcal{S}_c$. Next we describe the construction of any particular $\mathcal{S}_h$.

We fix $h$ and describe the construction of $\mathcal{S}_h$. Let $T_h = \{C_1, \ldots, C_k\}$, and let $P_h$ be a set of $k$ well-spread points in $\mathbb{R}^d$. Fix an arbitrary bijection $f_h : T_h \to P_h$, i.e. assign each clause to a distinct point from $P_h$. For notational convenience for a set of points $R \subseteq P_h$ we define $f_h^{-1}(R) = \{C_i \in T_h : f_h(C_i) \in R\}$. Let $\{H_1, \ldots, H_\ell\}$ be a maximal set of hyperplanes in $\mathbb{R}^d$ such that $\forall i \neq j : H_i^+ \cap P_h \neq H_j^+ \cap P_h$. For each $H_i$ we define a set $S_i \in \mathcal{S}_h$ to consist of the elements of the clauses in $f_h^{-1}(H_i^+ \cap P_h)$. We will occasionally abuse notation and also identify $S_i$ with the set of clauses $f_h^{-1}(H_i^+ \cap P_h)$. Note that this does not lead to ambiguity since the clauses in $T_h$ are pairwise disjoint. Finally, $\mathcal{S}_h = S_1 \cup \ldots \cup S_\ell$.

If the original instance of MAX-2-2-SET-SPLITTING is splittable then every set in the new set system has discrepancy zero. Assume that the original instance is unsplittable. By Lemma 3, for any assignment $\chi : X \to \{-1, 1\}$ there exists an independent set $T_h = \{C_1, \ldots, C_k\}$, s.t. at least $\Omega(N)$ of $C_1, \ldots, C_k$ are unsplit. Define, as before,

$$x_j^{(h)} = \sum_{e_i \in C_j} \chi(e_i). \tag{9}$$

We get an assignment of $P_h$ with values in the set $\{0, \pm 2, \pm 4\}$ by assigning each point $p_j \in P_h$ the value $x_j^{(h)}$ of its corresponding clause $C_j = f^{-1}(p_j)$. We focus on the points in $P_h$ with non-zero values. Consider the set $P_h' = \{f(C_j) : x_j^{(h)} \neq 0\}$. By Lemma 4, $P_h'$ is well-spread, and, therefore, by Lemma 5, there exists a hyperplane $H_i$ and a corresponding set $S_i \in \mathcal{S}_h$ s.t.

$$\sum_{j : C_j \in S_i} x_j^{(h)} = \Omega(|P_h'|^{1/2 - 1/2d}) = \Omega(N^{1/2 - 1/2d}). \tag{10}$$

It remains to show that $\pi_{\mathcal{S}}(s) = O(s^d)$ for $\mathcal{S} = \mathcal{S}_1 \cup \ldots \cup \mathcal{S}_c$. For a specific $\mathcal{S}_h$ and its corresponding set of points $P$ in $\mathbb{R}^d$, let $(P, \mathcal{H})$ be the set system consisting of all sets of the form

9

$H_i^+ \cap P$, where $H_i$ belongs to the set of maximal hyperplanes $\{H_1, \ldots, H_\ell\}$. Let $(\mathcal{E}_h, \mathcal{S}_h)$ be the set system defined above. Note that $\mathcal{E}_h$ is a subset of the $n$ elements of the original instance of MAX-2-2-SET-SPLITTING.

**Lemma 6** $\pi_{\mathcal{S}_h}(s) \leq \pi_{\mathcal{H}}(s)$.

**Proof:** Let $Y \subseteq \mathcal{E}_h$, and $|Y| = s$. Define $U = \{C_j \in T_h : Y \cap C_j \neq \emptyset\}$ to be the set of clauses intersected by $Y$, and $R = \{f(C_j) : C_j \in U\}$ to be the corresponding set of points in $P_h$. Since the clauses in $T_h$ are pairwise disjoint, $|R| = |U| \leq |Y| = s$.

We claim that $|(\mathcal{S}_h)_{|Y}| = |\mathcal{H}_{|R}|$. Notice that the claim completes the proof since $|\mathcal{H}_{|R}| \leq \pi_{\mathcal{H}}(|R|) \leq \pi_{\mathcal{H}}(s)$. To prove the claim we will establish a bijective map between sets $S_i \cap Y$ for $S_i \in \mathcal{S}_h$ and sets $H_i \cap R$. Associate each $S_i$ with its corresponding hyperplane $H_i$, i.e. $H_i = \{f(C_j) : C_j \in S_i\}$; then the bijective map sends a set $S_i \cap Y$ to $H_i \cap R$. To establish that this is a bijection we show that $S_i \cap Y \neq S_{i'} \cap Y$ if and only if $H_i \cap R \neq H_{i'} \cap R$. Indeed if $S_i \cap Y \neq S_{i'} \cap Y$, then there exists $a \in \mathcal{E}_h$ such that, without loss of generality, $a \in S_i \cap Y$ and $a \notin S_i' \cap Y$ (or vice versa). Then define the sets $A = \{C_j \in S_i : C_j \cap Y \neq \emptyset\}$ and $B = \{C_j \in S_{i'} : C_j \cap Y \neq \emptyset\}$. Since $a$ belongs to $S_i \cap Y$ but not to $S_i' \cap Y$, there must be some $C_j \in A$ such that $C_j \notin B$, since each element appears in a unique clause. Then the sets $\{f(C_j) : C_j \in A\} = H_i \cap R$ and $\{f(C_j) : C_j \in B\} = H_{i'} \cap R$ are also different. For the other direction, note that if $H_i \cap R \neq H_{i'} \cap R$, then $f^{-1}(H_i \cap R) \neq f^{-1}(H_{i'} \cap R)$. Thus, without loss of generality, there is some clause $C_j \in f^{-1}(H_i \cap R)$ such that $C_j \notin f^{-1}(H_{i'} \cap R)$ (or vice versa). Since the point corresponding to clause $C_j$, is in $R$, (i.e. $f(C_j) \in R$) at least one of the elements $a \in C_j$ belongs to $Y$, and thus belongs to $S_i \cap Y$. However, since $C_j \notin f^{-1}(H_{i'} \cap R)$, the elements in $C_j$ do not belong to $S_{i'}$. Thus, $a \notin S_{i'} \cap Y$, and we conclude that $S_i \cap Y \neq S_{i'} \cap Y$.□

By a simple union bound, for any $Y$, $|Y| = s$,

$$|\mathcal{S}_{|Y}| \leq \sum_{h=1}^{c} |(\mathcal{S}_h)_{|Y}| \leq \sum_{h=1}^{c} \pi_{\mathcal{S}_h}(s). \tag{11}$$

Using (11), Lemma 6, and the fact $\pi_{\mathcal{H}}(s) = O(s^d)$, we can derive the final bound on the shatter function of $\mathcal{S}$:

$$\pi_{\mathcal{S}}(s) = O(cs^d) = O(s^d). \tag{12}$$

# 6 A General Approach

In proving Theorem 2 and Theorem 5, we use a common approach: we start with an instance of MAX-2-2-SET-SPLITTING and compose this instance with a high discrepancy set system. This construction is then used to prove hardness for the resulting class of set systems. In Theorem 2 the high discrepancy set system is a Hadamard set system; In Theorem 5, a system of halfspaces is a high discrepancy set system. In the latter setting, we prove that the class of set systems resulting from the composition preserves the property of bounded shatter function. Here, we formalize this general approach.

As we have done previously, let $A$ refer to both a set system and to its incidence matrix. Consider a class of set systems $\mathcal{W} = \{W_k\}$, where $k$ denotes the number of elements in the respective set system, and such that for all $k$ in some infinite set $\mathcal{I}$ of positive integers, such a $W_k$ exists. Recall that if an instance $C$ of MAX-2-2-SET-SPLITTING on $n$ elements and $m$ clauses is unsplittable, then each assignment results in at least $\phi m$ unsplit clauses. Furthermore, recall that each element appears in at most $b$ clauses, and we can decompose $C$ into $c$ sets of independent clauses, $\{T_1, T_2, \ldots T_c\}$, where $c = 4b + 1$. For $\gamma = \phi/c$, define a discrepancy function $D'_\infty(W_k)$ as follows:

$$D'_\infty(W_k) = \min_{\substack{x \in \{0, \pm 2, \pm 4\}^k \\ \|x\|_2^2 \geq \gamma k}} \|W_k x\|_\infty. \tag{13}$$

For matrices $A$ and $B$ with equal numbers of columns, define $A \cup B$ as the matrix consisting of the rows of $A$ followed by the rows of $B$. With respect to the set systems represented by the incidence matrices $A$ and $B$, $A \cup B$ represents the union of the respective set systems. We abuse notation by allowing $T_i$ to represent both a set of clauses as well as a corresponding matrix in which each row represents a clause (and therefore has four '1's), and each column represents the clauses in which an element appears (and therefore has a single '1'). Furthermore, for convenience of notation, we assume that each set $T_i$ contains the same number of clauses. If this is not the case, note that we can pad $T_i$ with "dummy" clauses, each composed of four unique elements that appear only in that clause and not in any other clause in $T_i$ or any $T_j$ for $j \neq i$. Note that in any optimal assignment to the original MAX-2-2-SET-SPLITTING instance augmented with these "dummy" clauses, a "dummy" clause always has value zero. Moreover, note that we add at most $cm = O(m)$ "dummy" clauses using $4cm$ additional elements.

Let $\mathcal{A}$ be a class of set systems that, for each $k \in \mathcal{I}$, contains the set systems that can be written as $W_k T_1 \cup \ldots \cup W_k T_c$, for some $W_k$ where $D'_\infty(W_k) = f(k)$. Here, each $W_k T_h$ is the conventional matrix product. Note that for $W_k$ a 0-1 matrix, and $T_h$ a 0-1 matrix with at most one '1' per column, $W_k T_h$ is also a 0-1 matrix, and, therefore, can be interpreted as an incidence matrix.

**Theorem 6** *Given a set system $A \in \mathcal{A}$ on $N$ elements, it is NP-hard to distinguish between the cases* **(1)** $D_\infty(A) = 0$, *and* **(2)** $D_\infty(A) = f(\Omega(N))$.

**Proof:** For any $y \in \{-1, 1\}^N$, we can write:

$$Ay = W_m T_1 y \cup \ldots \cup W_m T_c y. \tag{14}$$

Note that, by construction, $N = \Theta(m)$. Also, note that $T_h y \in \{0, \pm 2, \pm 4\}^m$. Therefore, $\|Ay\|_\infty = \max_h \|W_m T_h y\|_\infty$. If for some $y$ and some $h$, we have $\|T_h y\|_2^2 \geq \gamma m$, then $\|Ay\|_\infty \geq D'_\infty(W_m)$.

Thus, if the original instance $C$ of MAX-2-2-SET-SPLITTING was splittable, then there exists a $y \in \{-1, 1\}^N$ such that $T_h y = 0$ for all $h$ and therefore $D_\infty(A) = 0$. If $C$ is unsplittable, then, by Lemma 3, for any $y \in \{-1, 1\}^N$, there exists an $h$ such that $\|T_h y\|_2^2 \geq \gamma m$. Therefore, $D_\infty(A) \geq D'_\infty(W_m) = f(m)$. Thus, we have $D_\infty(A) = f(\Omega(N))$. $\qquad\square$

# 7    Acknowledgments

We thank Nikhil Bansal for an enlightening conversation about the DISCREPANCY problem. We thank Bernard Chazelle for suggesting extensions of our main theorem; in particular to set systems of bounded VC-dimension/shatter function. We also thank him for making his useful book [Cha91] easily accessible and available for free on the Internet.

# References

[Ale90]    R. Alexander. Geometric methods in the study of irregularities of distribution. *Combinatorica*, 10(2):115–136, 1990.

[Ban10]    Nikhil Bansal.    Constructive algorithms for discrepancy minimization.    *CoRR*, abs/1002.2259, 2010.

[BS96]     J. Beck and V.T. Sós. Discrepancy theory. In *Handbook of combinatorics (vol. 2)*, page 1446. MIT Press, 1996.

[CGW05]   Moses Charikar, Venkatesan Guruswami, and Anthony Wirth. Clustering with qualitative information. *J. Comput. Syst. Sci.*, 71(3):360–383, 2005.

[Cha91]    Bernard Chazelle. *The Discrepancy Method.* Cambridge University Press, 1991.

[CMS95]   B. Chazelle, J. Matoušek, and M. Sharir. An elementary approach to lower bounds in geometric discrepancy. *Discrete and Computational Geometry*, 13(1):363–381, 1995.

[Gur03]    Venkatesan Guruswami. Inapproximability results for set splitting and satisfiability problems with no mixed clauses. *Algorithmica*, 38(3):451–469, 2003.

[Mat95]    J. Matoušek.  Tight upper bounds for the discrepancy of half-spaces.  *Discrete and Computational Geometry*, 13(1):593–601, 1995.

[Mat10]    J. Matousek. *Geometric Discrepancy: An Illustrated Guide.* Springer Verlag, 2010.

[Spe85]    Joel Spencer. Six standard deviations suffice. *Trans. Amer. Math. Soc.*, 289:679–706, 1985.

# A    The Generalized Alexander Lower Bound

The generalized Alexander lower bound (Lemma 5) follows easily from the proof technique of Chazelle et al. [CMS95]. We sketch how their proof needs to be modified. Our exposition follows section 3.3. of Chazelle [Cha91].

First, we introduce notation that closely follows Chazelle's. Let $X$ be a well-spread point set in $\mathbb{R}^d$, and let $\mathbf{v} = (v_1, 0, \dots, 0)$ be a vector in $\mathbb{R}^d$, where $v_1$ is a small real number to be specified

later. We consider a union of $X$ with $t = \lceil d/2 \rceil + 1$ copies of itself, each translated by a multiple of $\mathbf{v}$:

$$X_v = \bigcup_{j=0}^{t} (X + j\mathbf{v}). \tag{15}$$

Fix an assignment $\chi : X \to \{\pm 1, \pm 2\}$. The coloring is extended to $X_v$ as follows:

$$\chi(x + j\mathbf{v}) = (-1)^j \binom{t}{j} \chi(x). \tag{16}$$

Let $D(H)$ denote the discrepancy of $H^+ \cap X$, and let $D_v(H)$ denote the discrepancy of $H^+ \cap X_v$ with respect to the extended coloring. Consider a cube that encloses $X$, and pick a random hyperplane through the cube according to the motion-invariant measure on hyperplanes. By the probabilistic method, $\mathbb{E}[D(H)^2] \geq \max_H D(H)^2$, where the expectation is taken over picking a random hyperplane as described above. Chazelle [Cha91] shows that

$$\mathbb{E}[D(H)^2] = \Omega(\mathbb{E}[D_v(H)^2]). \tag{17}$$

The next step in the proof is to lower bound $\mathbb{E}[D_v(H)^2]$. Define a weight function $G(x, y)$ as

$$G(x, y) = \begin{cases} \sum_{j=-t}^{t} (-1)^j \binom{2t}{t+j} |x - y + j\mathbf{v}| & \text{if } x \neq y, \\ \\ -\binom{2t-2}{t-1} ||\mathbf{v}|| & \text{if } x = y. \end{cases} \tag{18}$$

Chazelle further proves the following facts:

$$\mathbb{E}[D_v(H)^2] = -\sum_{x,y \in X} \chi(x)\chi(y)G(x, y); \tag{19}$$

$$\sum_{x \neq y} |G(x, y)| = O(||\mathbf{v}||^{2t} n^{1+(2t-1)/d}). \tag{20}$$

All the statements so far are independent of the range of the assignment function $\chi$. Next we show how to modify the proof in order to accommodate the larger domain of assignments.

We separate the cross terms in the expression (19) for $\mathbb{E}[D_v(H)^2]$, and show that even if the points in $X$ are assigned colors from $\{\pm 1, \pm 2\}$, the cross terms are dominated by the remaining terms. Note that for any $x, y \in X$, $|\chi(x)\chi(y)| \leq 4$, and $\chi(x)^2 \geq 1$. Then,

$$\mathbb{E}[D_v(H)^2] = -\sum_{x} \chi(x)^2 G(x, x) - \sum_{x \neq y} \chi(x)\chi(y)G(x, y) \tag{21}$$

$$\geq -\sum_{x} G(x, x) - 4\sum_{x \neq y} |G(x, y)| \tag{22}$$

$$\tag{23}$$

By the definition of $G(x, x)$, and the bound (20), we have

$$\mathbb{E}[D_v(H)^2] = \Omega(n||\mathbf{v}|| - 4||v||^{2t} n^{1+(2t-1)/d}). \tag{24}$$

Setting $||v|| = cn^{-1/d}$ gives $\mathbb{E}[D_v(H)^2] = \Omega((c - 4c^{2t})n^{1-1/d})$. Choosing $c$ small enough so that $c - 4c^{2t} > 0$ completes the proof.

13