

## THE DISTANCE APPROACH TO APPROXIMATE COMBINATORIAL COUNTING

A. BARVINOK AND A. SAMORODNITSKY

### Abstract

We develop general methods to obtain fast (polynomial time) estimates of the cardinality of a combinatorially defined set via solving some randomly generated optimization problems on the set. Examples include enumeration of perfect matchings in a graph, linearly independent subsets of a set of vectors and colored spanning subgraphs of a graph. Geometrically, we estimate the cardinality of a subset of the Boolean cube via the average distance from a point in the cube to the subset with respect to some distance function. We derive asymptotically sharp cardinality bounds in the case of the Hamming distance and show that for small subsets a suitably defined “randomized” Hamming distance allows one to get tighter estimates of the cardinality.

### 1 Introduction

A general problem of combinatorial counting can be stated as follows: given a family  $\mathcal{F} \subset 2^X$  of subsets of the ground set  $X$ , compute or estimate the cardinality  $|\mathcal{F}|$  of the family. We would like to do the computation efficiently, *in polynomial time*. Of course, one should clarify what “given” means, especially since in most interesting cases  $|\mathcal{F}|$  is exponentially large in the cardinality  $|X|$  of the ground set. We assume that the family  $\mathcal{F}$  is defined by its *Optimization Oracle*:

#### 1.1 Optimization Oracle defining a family $\mathcal{F} \subset 2^X$ .

**Input:** A set of integer weights  $\gamma_x : x \in X$ .

**Output:** The number  $\min_{Y \in \mathcal{F}} \sum_{x \in Y} \gamma_x$ .

That is, for any given integer weighting  $\{\gamma_x\}$  on the set  $X$ , we should be able to produce the minimum weight of a subset  $Y \in \mathcal{F}$ . The following example was our main motivation.

---

The research of the first author was partially supported by NSF Grant DMS 9734138. The research of the second author was partially supported by a State of New Jersey grant.

EXAMPLE 1.2. PERFECT MATCHINGS IN A GRAPH. Let  $G = (V, E)$  be a graph with the set  $V$  of vertices and set  $E$  of edges. We assume that  $G$  has no loops (edges whose endpoints coincide). A set  $Y \subset E$  of edges is called a *matching* in  $G$  if every vertex of  $G$  is incident to at most one edge from  $Y$ . A matching  $Y$  is called *perfect* if every vertex of  $G$  is incident to precisely one edge from  $Y$ . Let  $\mathcal{F} \subset 2^E$  be the set of all *perfect matchings* in  $G$ . The problem of computing or estimating  $|\mathcal{F}|$  efficiently is one of the hardest and most intriguing problems of combinatorial counting, see, for example, [LoP], [JS1,2], [J] and [JSV]. It is known that the problem of exact counting of perfect matchings is hard. It belongs to the class of  $\#$  P-hard problems, see Chapter 18 of [P] for discussion of computational complexity in enumeration problems. Recently, Jerrum, Sinclair and Vigoda [JSV] found a way to estimate  $|\mathcal{F}|$  within any prescribed relative error  $\epsilon > 0$  in time polynomial in  $n$  and  $\epsilon^{-1}$  when  $G$  is a bipartite graph (in fact, they solved a more general problem of approximating the permanent of a non-negative matrix). However, for general graphs  $G$  the problem of efficient approximation of  $|\mathcal{F}|$  remains open.

We observe that Optimization Oracle 1.1 can be efficiently constructed. Indeed, if we assign integer weights  $\gamma_e: e \in E$  to the edges of the graph, the minimum weight of a perfect matching can be computed in  $O(|V|^3)$  time, see, for example, Section 11.3 of [PS].

The following example shows that sometimes optimization is nearly trivial but counting is still hard.

EXAMPLE 1.3. LINEARLY INDEPENDENT SUBSETS. Let  $\mathbb{F}$  be a field and let  $X \subset \mathbb{F}^d$  be a finite set of vectors. Given a number  $k \leq d$ , let  $\mathcal{F} \subset 2^X$  be the set of all linearly independent  $k$ -subsets of  $X$ . Optimization Oracle 1.1 for  $\mathcal{F}$  is supplied by the following *greedy algorithm*: given integer weights  $\{\gamma_x\}$  on  $X$ , we construct a linearly independent  $k$ -subset  $Y \subset X$  of the minimum weight by successively choosing vectors  $x_1, x_2, \dots, x_k$  of the minimum possible weight such that each set  $\{x_1\}, \{x_1, x_2\}, \dots, \{x_1, \dots, x_k\}$  is linearly independent. When  $\mathbb{F} = \mathbb{Q}$ , special cases of the counting problem for  $\mathcal{F}$  include counting *forests* (acyclic subgraphs) with  $k$  edges in a given graph and counting *spanning subgraphs* (connected subgraphs containing all the vertices of the graph) with  $k$  edges in a given connected graph. The last problem has interesting relations to percolation and network reliability. These counting problems appear to be difficult, except in some special cases, see [JS2] for a discussion, generalization to counting in *matroids* and some interesting conjectures and Section 12.4 of [PS] for optimization via

greedy algorithms.

Finally, we give an example of a more complicated structure, where no existing approaches to efficient counting seem to work.

**EXAMPLE 1.4. COLORED SPANNING SUBGRAPHS.** Let  $G = (V, E)$  be a connected graph with the set  $V$  of vertices and set  $E$  of edges. Suppose that the edges of  $G$  are colored in  $r$  colors (that is, the set  $E$  is represented as a union  $E = E_1 \cup \dots \cup E_r$  of  $r$  non-empty disjoint subsets). Let  $\mathcal{F} \subset 2^E$  consist of all subsets of edges such that the underlying graph is a spanning subgraph of  $G$  containing precisely one edge of each color. Optimization Oracle 1.1 for  $\mathcal{F}$  can be efficiently constructed using the matroid intersection algorithm (see, for example, Section 7.5 of [GLS]). The problem of counting or estimating  $|\mathcal{F}|$  appears to be rather difficult.

Next, we would like to discuss what “polynomial time” means. Typically, we are dealing with an infinite family of counting problems  $(\mathcal{F}_i, X_i) : i \in I$  (for example,  $X_i$  may range over the sets of edges of all finite graphs  $G_i$  and  $\mathcal{F}_i \subset 2^{X_i}$  may be the set of all perfect matchings in  $G_i$ ). We would like to construct an algorithm which works for every particular instance  $(\mathcal{F}_i, X_i)$ . If there exists a univariate polynomial *poly* such that the running time of the algorithm on the instance  $(\mathcal{F}_i, X_i)$  is bounded by  $\text{poly}(|X_i|)$ , we say that the algorithm is polynomial time. Our algorithms are *randomized*, that is, the algorithms rely on some coin tossing on the way, so the outcome is a random variable, which, with some high probability (say, 0.9) satisfies the desired properties. The probability of success can be made arbitrarily close to 1 by running the algorithm several times and taking a version of the majority vote (cf. Theorem 3.6 and remark that follows). For a general reference in the area of computational complexity and algorithms, see [P].

The most general approach to combinatorial counting has been via Monte Carlo method. The key component of the method is the ability to sample a random point from the (almost) uniform distribution on  $\mathcal{F}$ . Often, to achieve this, a Markov chain on the set  $\mathcal{F}$  is generated, so that it converges rapidly to the uniform distribution on  $\mathcal{F}$  (see [JS2] for a survey). Spectacular successes of this approach are finding a polynomial time randomized algorithm to count matchings of all sizes in a given graph [JS1] and to count perfect matchings in a given bipartite graph [JSV], both within any prescribed relative error. When the Markov chain approach works, it produces incomparably better results than the method of this paper. However, for many important counting problems, some of which are mentioned above, it is either not clear how to generate a rapidly mixing Markov chain

(as in Example 1.4) or, when there is a natural candidate (as in Example 1.3, see [JS2]), it seems to be extremely hard to prove that the chain is indeed converging rapidly enough to the steady state. In contrast, our approach produces very crude bounds, but it is totally insensitive to the fine structure of  $\mathcal{F}$ , so it is ready to handle a broad class of problems. A slight modification of the problem (for example, given small integer weights on the edges of the graph in Example 1.2, estimate the number of perfect matchings of the prescribed total weight) may lead to drastic changes in the construction of the underlying Markov chain but has almost no effect in our approach.

**1.5 The distance approach.** The main idea of our approach is as follows. Given a family  $\mathcal{F}$ , we identify it with a subset  $F$  of a metric space  $(\Omega, d)$ , such that for any given point  $x \in \Omega$  the distance  $d(x, F) = \min_{y \in F} d(x, y)$  can be quickly computed using Optimization Oracle 1.1 for  $\mathcal{F}$ . Then we estimate the cardinality  $|F|$  from the distance  $d(x, F)$  for a typical  $x \in \Omega$ . Intuitively, if  $|F|$  is small, we expect the distance  $d(x, F)$  from a random point  $x \in \Omega$  to be large and vice versa. In this paper,  $\Omega$  is the Boolean cube  $\{0, 1\}^n$  and  $d$  is either the Hamming distance or its modification, although as we discuss in section 5, some other possibilities may be of interest. Thus our approach can be considered as a refinement of the classical Monte Carlo method: we do not only register how often a randomly sampled point  $x \in \Omega$  lands in the target set  $F$ , but also take into account the distance  $d(x, F)$ . This allows us to get non-trivial bounds even when  $|F|$  is exponentially small with respect to  $|\Omega|$  so that  $x$  typically misses  $F$ .

The paper is organized as follows.

In section 2, we introduce a “geometric cousin” of Optimization Oracle 1.1. Distance Oracle 2.2 describes a subset  $F$  of the Boolean cube  $\{0, 1\}^n$  by computing a suitably defined distance  $d$  from a given point in the cube to the set. We show how to construct embeddings  $\phi : \mathcal{F} \rightarrow \{0, 1\}^n$ , so that the Distance Oracle for the image  $F = \phi(\mathcal{F})$  is derived from the Optimization Oracle for  $\mathcal{F}$ . We show that in some important cases (for example, when  $\mathcal{F}$  is the set of Example 1.2 of perfect matchings in a graph or the set of colored spanning subgraphs of Example 1.4), we can “squeeze”  $\mathcal{F}$  into a substantially smaller cube than we would have expected for a general family  $\mathcal{F}$ .

In section 3, we describe the bounds obtained by choosing  $d$  to be the Hamming distance in the cube. The bounds are sharp, meaning that we

can't possibly estimate (in polynomial time) the cardinality of a subset  $F \subset \{0, 1\}^n$  better if the only information available is the Hamming distance from any given point  $a \in \{0, 1\}^n$  to the set  $F$ .

In section 4, we describe how to get better bounds for small sets by using a suitably defined "randomized Hamming distance", which ignores a (random) part of the information contained in the standard Hamming distance. The isoperimetric problems arising here seem to be interesting in their own right. The proofs are not complicated but somewhat lengthy and therefore postponed till section 6.

In section 5, we discuss the types of estimates which can be obtained for particular counting problems (such as in Examples 1.2–1.4), possible ramifications of our approach and its relations with the Monte Carlo method.

## 2 Distance Oracle and Cubical Embeddings

The idea of our method is to represent  $\mathcal{F}$  geometrically as a subset  $F$  of the Boolean cube and then derive estimates of  $|\mathcal{F}|$  using the average distance from a point in the cube to  $F$ .

DEFINITIONS 2.1. Let  $C_n = \{0, 1\}^n$  be the Boolean cube and let  $\text{dist}$  be the Hamming distance in  $C_n$ , that is

$$\text{dist}(a, b) = \sum_{i:\alpha_i \neq \beta_i} 1 \quad \text{for } a = (\alpha_1, \dots, \alpha_n), b = (\beta_1, \dots, \beta_n) \in C_n.$$

More generally, let us fix  $n$  functions  $d_i : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{Z}$ ,  $i = 1, \dots, n$ , which we interpret as *penalties*. We assume that  $d_i \geq 0$  and that  $d(0, 0) = d(1, 1) = 0$ . Finally, let

$$d(a, b) = \sum_{i=1}^n d_i(\alpha_i, \beta_i) \quad \text{where } a = (\alpha_1, \dots, \alpha_n) \quad \text{and} \quad b = (\beta_1, \dots, \beta_n)$$

be the  $L^1$  distance function determined by the penalties  $\{d_i\}$ . If  $d_i(\alpha, \beta) = 1$  whenever  $\alpha \neq \beta$  then  $d(a, b) = \text{dist}(a, b)$ .

For a subset  $B \subset C_n$  and a point  $a \in C_n$ , let

$$d(a, B) = \min_{b \in B} d(a, b)$$

be the distance from  $a$  to  $B$ . In particular, let

$$\text{dist}(a, B) = \min_{b \in B} \text{dist}(a, b)$$

be the Hamming distance from a point  $a$  to the subset  $B$ .

We will be working with the following "geometric cousin" of Optimization Oracle 1.1.

## 2.2 Distance Oracle defining a set $F \subset C_n$ .

**Input:** A point  $a \in C_n$  and penalties  $d_i : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{Z}$ ,  $i = 1, \dots, n$ .

**Output:** The number  $d(a, F)$ .

There is an obvious way to associate a subset  $F \subset C_n$  with a family  $\mathcal{F} \subset 2^X$ ,  $n = |X|$ , of the Boolean cube.

**2.3 Straightforward embedding.** Let us identify the ground set  $X$  with the set  $\{1, \dots, n\}$ ,  $n = |X|$ . Let  $\mathcal{F}$  be a family of subsets of  $\{1, \dots, n\}$  given by its Optimization Oracle. For a subset  $Y \in \mathcal{F}$  let us define the indicator  $y \in C_n$ ,  $y = (\eta_1, \dots, \eta_n)$  by

$$\eta_i = \begin{cases} 1 & \text{if } i \in Y, \\ 0 & \text{if } i \notin Y. \end{cases}$$

Let  $F = \{y \in C_n : Y \in \mathcal{F}\}$  be the set of all indicators of subsets  $Y \in \mathcal{F}$ .

Let us construct the Distance Oracle for the set  $F \subset C_n$ . Given a point  $a = (\alpha_1, \dots, \alpha_n) \in C_n$  and penalties  $d_i$ ,  $i = 1, \dots, n$ , let us define weights  $\gamma_i$  by  $\gamma_i = d_i(\alpha_i, 1) - d_i(\alpha_i, 0)$ . Then for a set  $Y \subset \{1, \dots, n\}$  and its indicator  $y = (\eta_1, \dots, \eta_n) \in C_n$ , we have

$$\begin{aligned} \sum_{i \in Y} \gamma_i &= \sum_{i \in Y} (d_i(\alpha_i, 1) - d_i(\alpha_i, 0)) \\ &= \sum_{i=1}^n d_i(\alpha_i, \eta_i) - \sum_{i=1}^n d_i(\alpha_i, 0) = d(a, y) - d(a, 0). \end{aligned}$$

Hence, given the output

$$\lambda = \min_{Y \in \mathcal{F}} \sum_{i \in Y} \gamma_i$$

of Oracle 1.1 for the family  $\mathcal{F}$ , we can easily compute the output

$$d(a, F) = \lambda + d(a, 0)$$

of Oracle 2.2 for the set  $F$ . Thus, given an Optimization Oracle 1.1 for a family  $\mathcal{F} \subset 2^X$ , we can efficiently construct a Distance Oracle 2.2 for a set  $F \subset C_n$ ,  $n = |X|$ , such that  $|F| = |\mathcal{F}|$ .

To be able to estimate the cardinality  $|\mathcal{F}|$  with better precision, we would like to embed  $\mathcal{F}$  into a smaller Boolean cube. Sometimes this is indeed possible.

**2.4 Economical embedding.** Suppose that the ground set  $X$  can be represented as a union  $X = X_1 \cup \dots \cup X_k$  of (not necessarily disjoint) parts  $X_i$ , so that  $|Y \cap X_i| = 1$  for every subset  $Y \in \mathcal{F}$  and every  $X_i$ . In other words, every member of  $\mathcal{F}$  is a transversal of the cover of  $X$  by

$X_1, \dots, X_k$ . Let

$$m_i = \lceil \log_2 |X_i| \rceil \quad \text{and} \quad m = \sum_{i=1}^k m_i.$$

We construct an embedding  $\mathcal{F} \rightarrow C_m$  as follows.

First, we index the elements of  $X_i$  by distinct binary strings of length  $m_i$ , that is, we choose an embedding  $\phi_i : X_i \rightarrow C_{m_i}$ . Thus for any  $x \in X_i$  the point  $\phi_i(x)$  is a binary string of length  $m_i$  and  $\phi_i(x) \neq \phi_i(y)$  provided  $x \neq y$ .

Let us identify

$$C_m = C_{m_1} \times \dots \times C_{m_k}.$$

For a subset  $Y \in \mathcal{F}$ , let us define  $y \in C_m$  as

$$y = (y_1, \dots, y_k), \quad \text{where} \quad y_i = \phi_i(Y \cap X_i) \in C_{m_i}.$$

Note that  $y$  is well defined, since every intersection  $Y \cap X_i$  consists of a single point. Let  $F = \{y \in C_m : Y \in \mathcal{F}\}$ . Clearly,  $|F| = |\mathcal{F}|$ .

Given an Optimization Oracle 1.1 for  $\mathcal{F}$ , let us construct a Distance Oracle 2.2 for  $F$ . The input of Oracle 2.2 consists of a point  $a \in C_m$  (binary string of length  $m$ ) and penalty functions  $\{d_i : i = 1, \dots, m\}$ . We view  $a$  as

$$a = (a_1, \dots, a_k), \quad \text{where} \quad a_i \in C_{m_i}.$$

The penalties  $d_i, i = 1, \dots, m$  give rise to the  $L^1$  distance function  $d$  on binary strings, cf. Definition 2.1. For a point  $x \in X$ , let us define its weight  $\gamma_x$  by

$$\gamma_x = \sum_{i: x \in X_i} d(a_i, \phi_i(x)). \tag{2.4.1}$$

Let  $Y \in \mathcal{F}$  be a set and let  $y \in C_m$  be the point representing  $Y$ . We observe that

$$\sum_{x \in Y} \gamma_x = \sum_{x \in Y} \sum_{i: x \in X_i} d(a_i, \phi_i(x)) = \sum_{i=1}^k d(a_i, y_i) = d(a, y).$$

Hence, the outputs of Oracles 1.1 and 2.2 coincide:

$$\min_{Y \in \mathcal{F}} \sum_{x \in Y} \gamma_x = \min_{y \in Y} d(a, y).$$

Thus, given an Optimization Oracle 1.1 for a family  $\mathcal{F} \subset 2^X$ , we can efficiently construct a Distance Oracle 2.2 for a set  $F \subset C_m$ , such that  $|F| = |\mathcal{F}|$ . More precisely, given a point  $a \in C_m$  and penalties  $\{d_i\}$ , we compute weights  $\{\gamma_x\}$  on  $X$  by (2.4.1) in  $O(k|X| \ln |X|)$  time and then apply Optimization Oracle 1.1 to find the minimum weight  $\lambda$  of a subset  $Y \in \mathcal{F}$  in this weighting. The distance  $d(a, F)$  is equal to  $\lambda$ .

EXAMPLE 2.5. EMBEDDING PERFECT MATCHINGS AND COLORED SPANNING SUBGRAPHS. Let  $\mathcal{F}$  be the family of all perfect matchings in a graph  $G = (V, E)$ , see Example 1.2. The straightforward embedding 2.3 identifies  $\mathcal{F}$  with a subset  $F$  of the Boolean cube  $\{0, 1\}^{|E|}$  and provides us with Distance Oracle 2.2 for  $F$ . We will be better off using the economical embedding 2.4. Indeed, for a vertex  $v \in V$  of  $G$ , let  $E_v$  be the set of edges of  $G$  incident with  $v$ . Then  $E = \bigcup_{v \in V} E_v$  and every perfect matching has exactly one edge in every set  $E_v$ . Hence the embedding 2.4 identifies  $\mathcal{F}$  with a subset  $F$  of the Boolean cube  $\{0, 1\}^m$ , where

$$m = \sum_{v \in V} \lceil \log_2 |E_v| \rceil$$

and provides us with Distance Oracle 2.2 for  $F$ . Given a point  $a \in C_m$ , by (2.4.1) we compute weights  $\gamma_e$  on the edges  $E$  in  $O(|E| \ln |E|)$  time (since every edge  $e \in E$  belongs to exactly two sets  $E_v$ ) and then find the minimum weight  $\lambda$  of a perfect matching in  $G$  in  $O(|V|^3)$  time. The distance  $d(a, F)$  from  $a$  to  $F$  is equal to  $\lambda$ . Typically, if the graph has  $|V| = n$  vertices and  $\Omega(n^2)$  edges, the dimension of the straightforward embedding will be  $\Omega(n^2)$ , whereas the dimension of the economical embedding will be  $O(n \ln n)$ .

Similarly, if  $\mathcal{F}$  is the set of properly colored spanning subgraphs of Example 1.4, the partition  $E = E_1 \cup \dots \cup E_r$  gives rise to the economical embedding of  $\mathcal{F}$  into the Boolean cube  $\{0, 1\}^m$  of dimension  $m = \sum_{i=1}^r \lceil \log_2 |E_i| \rceil$ , as opposed to the dimension  $|E| = \sum_{i=1}^r |E_i|$  of the straightforward embedding.

### 3 Estimating Cardinality from the Hamming Distance

In this section, we obtain estimates of the cardinality of a subset  $F \subset C_n$  if we choose  $d_i(0, 1) = d_i(1, 0) = 1$ ,  $i = 1, \dots, n$  in Distance Oracle 2.2. In other words, we estimate  $|F|$ , provided we can compute the Hamming distance  $\text{dist}(x, F)$  to  $F$  from any given point  $x \in C_n$ , cf. Definitions 2.1. Our main tool is the *average* Hamming distance from a point to the set.

DEFINITION 3.1. Let  $A \subset C_n$  be a subset of the Boolean cube. Let

$$\Delta(A) = \frac{1}{2^n} \sum_{x \in C_n} \text{dist}(x, A)$$

be the average Hamming distance from a point in the cube to the set  $A$ .

Obviously,  $\Delta(A) \leq \Delta(B)$  if  $B \subset A$ .

EXAMPLE 3.2. SET CONSISTING OF A SINGLE POINT. Suppose that the set  $A$  is a point. Without loss of generality we assume that  $A = \{(0, \dots, 0)\}$ .



Then, for  $x = (\xi_1, \dots, \xi_n)$  we have  $\text{dist}(x, A) = \text{dist}(x, 0) = \xi_1 + \dots + \xi_n$  and

$$\Delta(A) = \frac{1}{2^n} \sum_{x \in C_n} \text{dist}(x, A) = \frac{1}{2^n} \sum_{x \in C_n} (\xi_1 + \dots + \xi_n) = \frac{n}{2}.$$

It follows then that  $\Delta(A) \leq n/2$  for any non-empty  $A \subset C_n$  and that  $\Delta(A) = n/2$  if and only if  $A$  consists of a single point.

Our first objective is to present a probabilistic algorithm that computes  $\Delta(A)$  approximately by averaging  $\text{dist}(x, A)$  for a number of randomly chosen  $x \in C_n$ .

**3.3 Algorithm for computing  $\Delta(A)$ .**

**Input:** A set  $A \subset C_n$  defined by its Distance Oracle 2.2 and a number  $\epsilon > 0$ .

**Output:** A number  $\alpha$  approximating  $\Delta(A)$  within error  $\epsilon$ .

**Algorithm:** Let  $k = \lceil 3n/2\epsilon^2 \rceil$ . Sample  $k$  points  $x_1, \dots, x_k \in C_n$  independently at random from the uniform distribution in the cube  $C_n$ . Apply Distance Oracle 2.2 to find  $\text{dist}(x_i, A)$ ,  $i = 1, \dots, k$ . Compute  $\alpha = \frac{1}{k} \sum_{i=1}^k \text{dist}(x_i, A)$ . Output  $\alpha$ .

To prove that Algorithm 3.3 indeed approximates  $\Delta(A)$  with the desired accuracy, we need a couple of technical results. The first lemma supplies us with important *concentration inequalities* for the Boolean cube.

LEMMA 3.4. *Let  $\xi_1, \dots, \xi_N$  be independent random variables taking values in  $\{0, 1\}$ . Let  $f : C_N \rightarrow \mathbb{R}$  be a function such that  $|f(x) - f(y)| \leq 1$  whenever  $\text{dist}(x, y) \leq 1$  and let  $\eta$  be the random variable  $f(\xi_1, \dots, \xi_N)$ . Then for any  $\delta > 0$*

$$\mathbf{P} \left\{ \eta : |\eta - \mathbf{E}(\eta)| \geq \delta \right\} \leq 2 \exp \left\{ -\frac{2\delta^2}{N} \right\}.$$

*Proof.* This is a special case of Lemma 1.2 of [M] . □

The next lemma provides a useful “scaling” trick.

LEMMA 3.5. *Let us fix positive integers  $k$  and  $n$  and let  $N = kn$ . Let us identify  $C_N = C_n \times \dots \times C_n = (C_n)^k$ . Thus a point  $x \in C_N$  is identified with a  $k$ -tuple  $x = (x_1, \dots, x_k)$ , where  $x_i \in C_n$  for  $i = 1, \dots, k$ .*

*For a subset  $A \subset C_n$ , let  $B = A \times \dots \times A = A^k \subset C_N$ . Then*

$$\text{dist}(x, B) = \sum_{i=1}^k \text{dist}(x_i, A) \quad \text{for any } x = (x_1, \dots, x_k) \in C_N$$

and

$$\Delta(B) = k\Delta(A).$$

*Proof.* Clearly,

$$\text{dist}(x, y) = \sum_{i=1}^k \text{dist}(x_i, y_i) \quad \text{for all } x, y \in C_N,$$

hence the first identity follows. Next,

$$\begin{aligned} \Delta(B) &= \frac{1}{2^N} \sum_{x \in C_N} \text{dist}(x, B) = \frac{1}{2^N} \sum_{x_1, \dots, x_k \in C_n} \sum_{i=1}^k \text{dist}(x_i, A) \\ &= \frac{k 2^{n(k-1)}}{2^{nk}} \sum_{x \in C_n} \text{dist}(x, A) = \frac{k}{2^n} \sum_{x \in C_n} \text{dist}(x, A) = k\Delta(A). \quad \square \end{aligned}$$

Now we can prove correctness of Algorithm 3.3.

**Theorem 3.6.** *With probability at least 0.9, the output  $\alpha$  of Algorithm 3.3 satisfies the inequality  $|\Delta(A) - \alpha| \leq \epsilon$ .*

*Proof.* Let  $N = nk$  and let us identify  $C_N = (C_n)^k$  as in Lemma 3.5. Let  $B = A^k \subset C_N$ . Let  $f : C_N \rightarrow \mathbb{R}$  be defined by  $f(x) = \text{dist}(x, B)$ . Interpreting  $x = (\xi_1, \dots, \xi_N)$  as a vector of  $N$  independent random variables  $\xi_i$  uniformly distributed on  $\{0, 1\}$ , applying Lemma 3.4 with  $\delta = k\epsilon$  and observing that  $\mathbf{E}(f) = \Delta(B)$ , we conclude that

$$\mathbf{P}\{x : |\text{dist}(x, B) - \Delta(B)| \geq k\epsilon\} \leq 2 \exp\left\{-\frac{2(\epsilon k)^2}{N}\right\} = 2 \exp\left\{-\frac{2\epsilon^2 k}{n}\right\} \leq 0.1.$$

Since by Lemma 3.5

$$\Delta(B) = k\Delta(A) \quad \text{and} \quad \frac{1}{k} \sum_{i=1}^k \text{dist}(x_i, A) = \frac{1}{k} \text{dist}(x, B)$$

for  $x = (x_1, \dots, x_k)$ , we conclude that

$$\begin{aligned} \mathbf{P}\left\{x_1, \dots, x_k : \left|\frac{1}{k} \sum_{i=1}^k \text{dist}(x_i, A) - \Delta(A)\right| \geq \epsilon\right\} \\ = \mathbf{P}\{x : |\text{dist}(x, B) - \Delta(B)| \geq k\epsilon\} \leq 0.1, \end{aligned}$$

and the proof follows. □

**REMARK.** Hence to evaluate  $\Delta(A)$  within error  $\epsilon$  we have to average  $O(n\epsilon^{-2})$  values  $\text{dist}(x_i, A)$ . By doing that, we allow probability 0.1 of failure. As usual, to attain a lower probability  $\delta > 0$  of failure, one should run Algorithm 3.3  $O(\ln \delta^{-1})$  times and then select the median of the computed  $\alpha$ 's (cf. [JVV]). For all applications, choosing  $\epsilon = 1$  will suffice and in many cases  $\epsilon = \sqrt{n}$  will do (cf. section 5.1). Hence, often we will have to apply Oracle 2.2 only a constant number of times.

We would like to relate the value of  $\Delta(A)$  to the cardinality  $|A|$ .

**DEFINITION 3.7.** ENTROPY FUNCTION. For  $0 \leq x \leq 1/2$  let

$$H(x) = x \log_2 \frac{1}{x} + (1 - x) \log_2 \frac{1}{1 - x}.$$

We agree that  $H(0) = 0$ . Thus  $H$  is an increasing concave function on the interval  $[0, 1/2]$ .

We use the following estimate (see, for example, Theorem 1.4.5 of [Li])

$$\sum_{k=0}^r \binom{n}{k} \leq 2^{nH(r/n)} \quad \text{for } r \leq n/2. \tag{3.7.1}$$

Also, we remark that around  $x = +0$  we have

$$H(x) = x \log_2 \frac{1}{x} + O(x) \quad \text{and} \quad H\left(\frac{1}{2} - x\right) = 1 - \frac{2}{\ln 2} x^2 + O(x^3) \tag{3.7.2}$$

We will use the classical isoperimetric inequality for the Boolean cube (see, for example, [L]).

**Theorem 3.8** (Harper’s theorem). *Let  $A \subset C_n$  be a set such that*

$$|A| \geq \sum_{k=0}^r \binom{n}{k}$$

*for some integer  $r$ . Then, for any non-negative integer  $t$*

$$|\{x \in C_n : \text{dist}(x, A) \leq t\}| \geq \sum_{k=0}^{r+t} \binom{n}{k}.$$

We are going to obtain an estimate of the cardinality of a set  $A \subset C_n$  in terms of the average Hamming distance  $\Delta(A)$  from a point  $x \in C_n$  to  $A$ . It is convenient to express the estimate in terms of a related quantity

$$\rho = \rho(A) = \frac{1}{2} - \frac{\Delta(A)}{n}.$$

As follows from Example 3.2, for every non-empty set  $A \subset C_n$  we have  $0 \leq \rho(A) \leq 1/2$ . We observe that  $\rho(A) = 0$  if and only if  $A$  consists of a single point and that  $\rho(A) = 1/2$  if and only if  $A$  is the whole cube  $C_n$ .

**Theorem 3.9.** *Let  $A \subset C_n$  be a non-empty set. Let*

$$\rho = \frac{1}{2} - \frac{\Delta(A)}{n}.$$

*Then*

$$1 - H\left(\frac{1}{2} - \rho\right) \leq \frac{\log_2 |A|}{n} \leq H(\rho).$$

Before we proceed with a formal proof, we would like to highlight some ideas.

**3.10 The idea of the proof. Extremal sets.** Let  $A \subset C_n$  be a set. Concentration inequalities (Lemma 3.4) imply that the average distance  $\Delta(A)$  is approximately equal to the distance  $\text{dist}(x, A)$  from a “typical” point  $x \in C_n$  to  $A$  (see also Sections 6.2 and 7.9 of [MiS] ). For a given positive integer  $t$ , let us consider the  $t$ -neighborhood  $A_t = \{x \in C_n : \text{dist}(x, A) \leq t\}$  of  $A$ . We expect that  $\Delta(A) \approx t_1$ , where  $t_1$  is the smallest value of  $t$  such that  $A_t$  covers “almost all” cube  $C_n$ . The neighborhood  $A_t$  grows the slowest when  $A$  is a ball in the Hamming metric, that is when  $A = \{x : \text{dist}(x, x_0) \leq r\}$  for some  $x_0 \in C_n$  and some  $r > 0$ , as follows from Harper’s theorem 3.8, cf. also [L] . Hence the upper bound for  $n^{-1} \log_2 |A|$  in Theorem 3.9 is attained (up to an  $O(n^{-1/2})$  error term) when  $A$  is a ball. The neighborhood  $A_t$  grows the fastest when the points of  $A$  are spread around in  $C_n$ . In any case, the size  $|A_t|$  does not exceed the sum of sizes of the balls of radius  $t$  centered at the points of  $A$ . Thus the lower bound for  $n^{-1} \log_2 |A|$  in Theorem 3.9 is obtained from this “packing” type argument. One can show that if the points of  $A$  are chosen at random in  $C_n$ , then with high probability the lower bound is indeed attained asymptotically. More precisely, let us fix a number  $0 < \beta < 1$  and let  $A$  be the set of  $\lfloor 2^{\beta n} \rfloor$  points chosen at random from  $C_n$ . Then with the probability that tends to 1 as  $n$  grows to infinity,  $\beta = 1 - H(\frac{1}{2} - \rho) + O(n^{-1/2})$ . The proof is straightforward, but technical and therefore omitted.

Finally, we note that using average distance  $\Delta(A)$  and the scaling trick (Lemma 3.5) allows us to get rid of  $O(n^{-1/2})$  error terms in the proof.

*Proof of Theorem 3.9.* Let us choose a positive even integer  $m$ , let  $N = mn$  and let us identify  $C_N = (C_n)^m$ , as in Lemma 3.5. Let  $B = A^m \subset C_N$ . Let us fix the uniform probability measure  $\mathbf{P}$  on  $C_N$ .

Let  $\alpha = \log_2 |A|/n$ , so  $|A| = 2^{\alpha n}$  and  $|B| = 2^{\alpha N}$ . Let  $0 \leq \gamma \leq 1/2$  be a number such that  $H(\gamma) = \alpha$  and let  $r = \lfloor N\gamma \rfloor$ . Then by (3.7.1)

$$|B| = 2^{N \cdot H(\gamma)} \geq \sum_{k=0}^r \binom{N}{k}.$$

Then Theorem 3.8 implies that

$$|\{x \in C_N : \text{dist}(x, B) \leq N/2 - r\}| \geq \sum_{k=0}^{N/2} \binom{N}{k} = 2^{N-1}.$$

Therefore,

$$\mathbf{P} \{x \in C_N : \text{dist}(x, B) \leq \frac{N}{2} - r\} \geq \frac{1}{2}.$$

We have that  $x = (x_1, \dots, x_m)$  for some  $x_i \in C_n$  and that  $\text{dist}(x, B) =$

$\text{dist}(x_1, A) + \dots + \text{dist}(x_m, A)$  (see Lemma 3.5). Therefore,

$$\mathbf{P} \left\{ (x_1, \dots, x_m) : \frac{1}{m} \sum_{i=1}^m \text{dist}(x_i, A) \leq \frac{N}{2m} - \frac{r}{m} \right\} \geq \frac{1}{2}. \tag{1}$$

Now we observe that

$$\frac{N}{2m} - \frac{r}{m} \longrightarrow \frac{n}{2} - n\gamma \quad \text{as } m \rightarrow +\infty. \tag{2}$$

Furthermore, by the Law of Large Numbers,

$$\frac{1}{m} \sum_{i=1}^m \text{dist}(x_i, A) \longrightarrow \Delta(A) \quad \text{in probability as } m \rightarrow +\infty. \tag{3}$$

Hence the assumption that  $\Delta(A) > n/2 - n\gamma$  would contradict (1)–(3). Thus we must have  $\Delta(A) \leq n/2 - n\gamma$ , which implies that  $\gamma \leq \rho(A)$ . Hence  $\alpha = H(\gamma) \leq H(\rho)$  and the upper bound is proven.

Let us prove the lower bound. We observe that for every point  $b \in C_N$  and any  $N/2 \geq s \geq 0$

$$|\{x \in C_N : \text{dist}(x, b) \leq s\}| = \sum_{k=0}^s \binom{N}{k} \leq 2^{N \cdot H(s/N)}.$$

Therefore,

$$|\{x \in C_N : \text{dist}(x, B) \leq s\}| \leq |B| 2^{N \cdot H(s/N)} = 2^{N \cdot (H(s/N) + \alpha)}.$$

Hence

$$\mathbf{P} \{x \in C_N : \text{dist}(x, B) \leq s\} \leq 2^{N \cdot (H(s/N) + \alpha - 1)}.$$

Therefore,

$$\mathbf{P} \left\{ (x_1, \dots, x_m) : \frac{1}{m} \sum_{i=1}^m \text{dist}(x_i, A) \leq s/m \right\} \leq 2^{N \cdot (H(s/N) + \alpha - 1)}. \tag{4}$$

If  $\Delta(A) = n/2$  then  $A$  is a point and the lower bound in Theorem 3.9 is satisfied. Otherwise, let us fix an  $\epsilon > 0$  such that  $(1 + \epsilon)\Delta(A)/n < 1/2$  and let  $s = \lceil m(1 + \epsilon)\Delta(A) \rceil$ . We have

$$s/m \longrightarrow (1 + \epsilon)\Delta(A) \quad \text{and} \quad s/N \longrightarrow (1 + \epsilon)\Delta(A)/n \quad \text{as } m \rightarrow +\infty. \tag{5}$$

Hence the assumption that  $H((1 + \epsilon)\Delta(A)/n) + \alpha - 1 < 0$  would contradict (3)–(5). Therefore,  $H((1 + \epsilon)\Delta(A)/n) + \alpha - 1 \geq 0$  for any  $\epsilon > 0$  and  $H(\Delta(A)/n) + \alpha - 1 \geq 0$ . Since  $\Delta(A)/n = \frac{1}{2} - \rho$ , the proof follows.  $\square$

For applications, the most interesting case is when  $n^{-1} \log_2 |A|$  is small, that is  $\rho \approx 0$ .

**COROLLARY 3.11.** *There exist positive constants  $c_1$  and  $c_2$  such that for any non-empty set  $A \subset C_n$  and for  $\rho = \frac{1}{2} - \frac{\Delta(A)}{n}$  we have*

$$c_1 \cdot \rho^2 \leq \frac{\ln |A|}{n} \leq c_2 \cdot \rho \ln \frac{1}{\rho}.$$

In particular, for any  $c_1 < 2$  and any  $c_2 > 1$ , the inequality holds in a sufficiently small neighborhood of  $\rho = 0$ .

*Proof.* Follows from Theorem 3.9 by (3.7.2).  $\square$

**3.12 Discussion.** Figure 1 depicts the feasible region for  $n^{-1} \log_2 |A|$  as described by Theorem 3.9. Thus possible values of  $n^{-1} \log_2 |A|$  with the

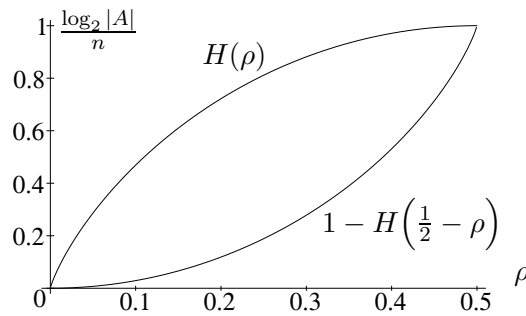


Figure 1

given value of  $\rho$  form a vertical interval between the two curves. As we discussed in section 3.10, asymptotically both bounds are sharp. Remarkably, the bounds converge at  $\rho = 0$  and  $\rho = 0.5$ . On the other hand, the difference is the greatest when  $\rho = 1/4$ . Thus, if the average Hamming distance from a point  $x \in C_n$  to a set  $A \subset C_n$  is  $n/4$ , the set  $A$  can contain as many as  $2^{0.811n}$  points and as few as  $2^{0.189n}$  points. We note that if  $A$  is a face (subcube) of the Boolean cube then the corresponding point lies on the straight line  $n^{-1} \log_2 |A| = 2\rho(A)$ .

Corollary 3.11 (with somewhat weaker constants and stated in different terms) together with the observation that the distance  $\text{dist}(x, A)$  for a randomly chosen point  $x \in C_n$  allows one to estimate  $\rho$  up to an  $O(n^{-1/2})$  error constitute the main result of the earlier paper [B]. Consequently, the main conclusion of [B] is equivalent to stating that the Hamming distance to  $A$  from a random point  $x$  in the Boolean cube allows one to decide whether  $|A|$  is exponentially large in  $n$ . In this paper, we make improvements in several directions. The most important one is that the optimization functionals of [B] are recognized as *distances*, which allows us to establish connections

with isoperimetric inequalities. In particular, Theorem 3.9 provides us with sharp bounds valid for all  $0 \leq \rho \leq 1/2$ . Also, using Algorithm 3.3 and Theorem 3.6 we get rid of the  $O(n^{-1/2})$  error term. This allows us to obtain meaningful cardinality estimates for sets with small values of  $\rho$ . Curiously, we can even distinguish in polynomial time between a set consisting of a single point ( $\rho = 0$ ) and a set having more than one point (one can show that  $\rho = \Omega(1/n)$  in that case), although apparently we can't distinguish between sets consisting of 2 and 3 points respectively. Finally, construction 2.4 of "economical embedding" allows us to obtain tighter bounds for a wide class of problems by lowering the dimension of the ambient Boolean cube.

In section 4, we show how using "randomized Hamming distance" allows one to get better estimates for sets  $A$  with  $n^{-1} \log_2 |A|$  small, which is the case in most applications. Essentially, the randomized Hamming distance will allow us to "sandwich" such a set between a random set and a face (subcube), as opposed to a random set and a Hamming ball in the case of the standard Hamming metric, cf. section 3.10.

#### 4 Randomized Hamming Distance

Let us fix a number  $0 < p \leq 1$  and let  $q = 1 - p$ . In this section, we construct a quantity  $\Delta(A, p)$ , which measures the cardinality of "small" subsets  $A \subset C_n$  of the Boolean cube in a somewhat more precise way than the average Hamming distance  $\Delta(A)$  discussed in section 3. In fact,  $\Delta(A, 1) = \Delta(A)$ , so  $\Delta(A)$  is a particular case of  $\Delta(A, p)$ .

DEFINITIONS 4.1. Let  $\Lambda_n$  be a copy of the Boolean cube  $\{0, 1\}^n$ . We make  $\Lambda_n$  a probability space by letting

$$\mathbf{P} \{l\} = p^{|l|} q^{n-|l|}, \text{ where } |l| = \lambda_1 + \dots + \lambda_n \text{ for } l = (\lambda_1, \dots, \lambda_n).$$

Hence a vector  $l = (\lambda_1, \dots, \lambda_n)$  from  $\Lambda_n$  is interpreted as a realization of  $n$  independent random variables  $\lambda_i$  such that  $\mathbf{P} \{\lambda_i = 1\} = p$  and  $\mathbf{P} \{\lambda_i = 0\} = q$ .

For  $x, y \in C_n$  and an  $l \in \Lambda_n$ , where  $x = (\xi_1, \dots, \xi_n)$ ,  $y = (\eta_1, \dots, \eta_n)$  and  $l = (\lambda_1, \dots, \lambda_n)$ , let

$$d_l(x, y) = \sum_{i:\xi_i \neq \eta_i} \lambda_i.$$

In other words, we count disagreement in the  $i$ -th coordinate of  $x$  and  $y$  if and only if the value of  $\lambda_i$  is 1. Thus if  $l = (1, \dots, 1)$ , we have  $d_l(x, y) = \text{dist}(x, y)$ , the usual Hamming distance.

For  $l \in \Lambda_n$  and a set  $A \subset C_n$ , let

$$d_l(x, A) = \min_{y \in A} d_l(x, y).$$

Finally, let

$$\Delta(A, p) = \sum_{l \in \Lambda_n} \sum_{x \in C_n} d_l(x, A) \frac{p^{|l|} q^{n-|l|}}{2^n}.$$

In other words,  $\Delta(A, p)$  is the expected value of  $d_l(x, A)$ , where  $x = (\xi_1, \dots, \xi_n)$  and  $l = (\lambda_1, \dots, \lambda_n)$  are vectors of independent random variables such that

$\mathbf{P} \{\lambda_i = 1\} = p$ ,  $\mathbf{P} \{\lambda_i = 0\} = q$  and  $\mathbf{P} \{\xi_i = 0\} = \mathbf{P} \{\xi_i = 1\} = 1/2$ . Obviously,  $\Delta(A, p) \leq \Delta(B, p)$  if  $B \subset A$ .

It follows that for a fixed non-empty  $A \subset C_n$ , the value  $\Delta(A, p)$  is a polynomial in  $p$  of degree at most  $n$ .

**EXAMPLE 4.2. SET CONSISTING OF A SINGLE POINT.** Suppose that the set  $A$  consists of a single point. Without loss of generality we assume that  $A = \{(0, \dots, 0)\}$ . Then for  $x = (\xi_1, \dots, \xi_n)$  and  $l = (\lambda_1, \dots, \lambda_n)$ ,

$$d_l(x, A) = \sum_{i=1}^n \lambda_i \xi_i.$$

Interpreting  $\lambda_i$  and  $\xi_i$ ,  $i = 1, \dots, n$  as independent random variables such that  $\mathbf{P} \{\xi_i = 1\} = \mathbf{P} \{\xi_i = 0\} = 1/2$  and  $\mathbf{P} \{\lambda_i = 1\} = p$ ,  $\mathbf{P} \{\lambda_i = 0\} = q$ , we get

$$\Delta(A, p) = \mathbf{E} \sum_{i=1}^n \lambda_i \xi_i = \sum_{i=1}^n (\mathbf{E} \lambda_i)(\mathbf{E} \xi_i) = \frac{np}{2}.$$

It follows then that for any non-empty set  $A \subset C_n$  we have  $\Delta(A, p) \leq np/2$  and that  $\Delta(A, p) = np/2$  if and only if  $A$  consists of a single point (we agreed that  $p > 0$ ).

As was the case with  $\Delta(A)$ , the functional  $\Delta(A, p)$  can be easily computed by averaging. For a set  $A \subset C_n$  defined by its Distance Oracle 2.2 and any  $l = (\lambda_1, \dots, \lambda_n)$  the value of  $d_l(x, A)$  is computed by choosing the penalties  $d_i(0, 1) = d_i(1, 0) = 1$  when  $\lambda_i = 1$  and  $d_i = 0$  when  $\lambda_i = 0$ .

### 4.3 Algorithm for Computing $\Delta(A, p)$ .

**Input:** A set  $A \subset C_n$  given by its Distance Oracle 2.2, a number  $1 \geq p > 0$  and an  $\epsilon > 0$ .

**Output:** A number  $\alpha$  approximating  $\Delta(A, p)$  within error  $\epsilon$ .

**Algorithm:** Let  $k = \lceil 3n/\epsilon^2 \rceil$ . Sample  $k$  points  $x_1, \dots, x_k \in C_n$  independently at random from the uniform distribution in  $C_n$  and  $k$  points



$l_1, \dots, l_k \in \Lambda_n$  independently at random from the distribution in  $\Lambda_n$ . Apply Distance Oracle 2.2 to compute  $d_{l_i}(x_i, A)$ ,  $i = 1, \dots, k$ . Compute  $\alpha = \frac{1}{k} \sum_{i=1}^k \text{dist}_{l_i}(x_i, A)$ . Output  $\alpha$ .

**Theorem 4.4.** *With probability at least 0.9, the output  $\alpha$  of Algorithm 4.3 satisfies the inequality  $|\Delta(A, p) - \alpha| \leq \epsilon$ .*

We postpone the proof till section 6.

We are going to obtain estimates of the cardinality  $|A|$  of a set  $A \subset C_n$  in terms of the quantity  $\Delta(A, p)$ . As in section 3, it is convenient to work with a related quantity

$$\rho = \rho(A, p) = \frac{p}{2} - \frac{\Delta(A, p)}{n}.$$

From Definitions 4.1, for any non-empty  $A \subset C_n$ , the function  $\rho(A, p)$  is a polynomial in  $p$  of degree at most  $n$ . As follows from Example 4.2,  $0 \leq \rho \leq p/2$  for any non-empty set  $A \subset C_n$ . Our estimate will be useful for “small” sets  $A$  where  $n^{-1} \ln |A|$  is close to 0.

**Theorem 4.5.** *Let  $A \subset C_n$  be a non-empty set. Let*

$$\rho = \frac{p}{2} - \frac{\Delta(A, p)}{n}.$$

Then

$$\frac{\rho^2}{p} \leq \frac{\ln |A|}{n}. \tag{4.5.1}$$

Suppose that  $\rho \leq 1/4$  and that

$$p \geq \frac{\ln 2 + \ln(1 - 2\rho)}{\ln(1 - 2\rho) - \ln(2\rho)}. \tag{4.5.2}$$

Then

$$\frac{\ln |A|}{n} \leq 2\rho \ln \frac{1}{2\rho} + (1 - 2\rho) \ln \frac{1}{1 - 2\rho}. \tag{4.5.3}$$

We obtain the following counterpart of Corollary 3.11.

**COROLLARY 4.6.** *Let us choose any  $c_3 < 1/(\ln 2) \approx 1.44$  and any  $c_4 > 2$ . Then there exists a  $\delta > 0$  such that for any non-empty  $A \subset C_n$  with  $n^{-1} \ln |A| \leq \delta$  there exists a  $0 < p \leq 1$  such that for  $\rho = \frac{p}{2} - \frac{\Delta(A, p)}{n}$  one has*

$$c_3 \cdot \rho^2 \ln \frac{1}{\rho} \leq \frac{\ln |A|}{n} \leq c_4 \cdot \rho \ln \frac{1}{\rho}.$$

*Proof.* By (4.5.1), we have  $\rho \leq \sqrt{n^{-1} \ln |A|} \leq \sqrt{\delta}$ , so  $\rho(A, p)$  is small if  $\delta$  is small, no matter what  $p$  is. We observe that for small positive  $\rho$  the right-hand side of (4.5.2) is of the order  $(\ln 2) \ln^{-1}(1/\rho)$  and the right-hand side of (4.5.3) is of the order  $2\rho \ln(1/\rho)$ .

Given  $c_3 < (\ln 2)^{-1}$  and  $c_4 > 2$ , let us choose  $1/16 > \delta > 0$  in such a way that the right-hand side of (4.5.2) does not exceed  $(c_3)^{-1} \ln^{-1}(1/\rho)$  and the right-hand side of (4.5.3) does not exceed  $c_4 \rho \ln(1/\rho)$  for all  $0 < \rho < \sqrt{\delta} < 1/4$ .

We recall that  $|A| = 1$  if and only if  $\rho = 0$ , in which case the bounds of Corollary 4.6 are satisfied by default. If  $|A| > 1$ , one can see (cf. also section 4.7 below) that  $\rho(A, p) = \Omega(p/n)$ . Given a set  $A \subset C_n$ ,  $|A| > 1$ , let us choose the smallest  $p \geq 0$  that satisfies the inequality (4.5.2). Then  $0 < p < 1$  since the right-hand side of (4.5.2) is  $\Omega(\ln^{-1}(n/p))$  (and therefore goes to 0 slower than  $p$ ), and smaller than 1 for  $0 < \rho < 1/4$ . Since  $\rho(A, p)$  depends continuously on  $p$ , we must have equality in (4.5.2) (otherwise, we could have taken a smaller  $p$ ). Thus  $p \leq (c_3)^{-1} \ln^{-1}(1/\rho)$  and the proof follows by (4.5.1)–(4.5.3).  $\square$

**4.7 Extremal sets.** Let us fix a  $0 < p \leq 1$  and an  $\epsilon > 0$ . Then there exists an  $\alpha = \alpha(p, \epsilon) > 0$  with the following property: if  $A \subset C_n$  is a set of  $\lfloor 2^{\alpha n} \rfloor$  points randomly chosen from the Boolean cube, then with the probability that tends to 1 as  $n$  grows to infinity,  $n^{-1} \ln |A| < (2 + \epsilon)\rho^2/p$ . Hence for any  $p > 0$  the bound (4.5.1) is tight up to a constant factor for sufficiently small random sets. The proof is rather technical and therefore omitted.

One can show that the bound (4.5.3) is asymptotically tight on small faces of the cube  $C_n$ . More precisely, let us fix a  $\delta > 0$  (to be adjusted later), let  $m = \lfloor \delta n \rfloor$  and let  $A \subset C_n$  be an  $m$ -dimensional face of the Boolean cube

$$A = \{(\xi_1, \dots, \xi_n) : \xi_i = 0 \text{ for } i = m + 1, \dots, n\}.$$

Thus  $|A| = 2^m$ . Moreover, a computation similar to that of Example 4.2 shows that  $\rho(A, p) = pm/2n$ . Hence we have

$$\frac{\ln |A|}{n} = \frac{2 \ln 2}{p} \rho(A, p).$$

We observe that  $\rho(A, p) \leq \delta/2$ . Hence for any small  $\epsilon > 0$  one can find  $\delta = \delta(\epsilon) > 0$  such that there exists  $p$  satisfying (4.5.2) and such that  $p < (1 + \epsilon)(\ln 2) \ln^{-1}(1/\rho)$ . For such a  $p$ , we have

$$\frac{\ln |A|}{n} \geq \frac{2}{1 + \epsilon} \rho \ln \frac{1}{\rho},$$

so the bound (4.5.3) is indeed asymptotically tight for small sets.

Apparently, the sets  $A$  having the largest cardinality among all sets with the given value of  $\rho(A, p)$  evolve from the balls in the Hamming metric

for  $p = 1$  (see section 3.10) to faces at  $p \rightarrow 0$ . Since faces are packed somewhat less tightly than balls, we gain in Corollary 4.6 as compared to Corollary 3.11.

The proof of Theorem 4.5 is postponed till section 6.

**4.8 Discussion.** Corollary 4.6 implies that for small sets  $A$  by “tuning up”  $p$  we can get an additional logarithmic factor which brings the lower bound for  $n^{-1} \ln |A|$  a little closer to the upper bound compared to the bound of Corollary 3.11. Any  $p$  which is only slightly bigger than the bound (4.5.2) will do. Suppose, for example, that  $A \subset C_n$  is a set such that  $n^{-1} \ln |A| \sim n^{-\alpha}$  for some  $0 < \alpha < 1$ . By applying Algorithm 3.3 to approximate  $\Delta(A) = \Delta(A, 1)$  and Theorem 3.9 to interpret the results, the worst lower bound we can get for  $n^{-1} \ln |A|$  is  $\sim n^{-2\alpha} / \ln^2 n$  (this happens when  $A$  is a ball in the Hamming metric, but we think it is a “random set”, see section 3.10) and the worst upper bound we can get is  $\sim n^{-\alpha/2} \ln n$  (this happens when  $A$  is a “random set” but we think that it is a ball). Now, by (4.5.1) it follows that  $\rho(A, p) = O(n^{-\alpha/2})$  for any  $p$ . Then we can choose some  $p = O(\ln^{-1} n)$  that satisfies (4.5.2). Applying Algorithm 4.3 to approximate  $\Delta(A, p)$  and Theorem 4.5 to interpret the results, for  $n^{-1} \ln |A|$  we would obtain a lower bound of the form  $\sim n^{-2\alpha} / \ln n$  at worst (this happens when  $A$  is a face but we think it is a random set) and an upper bound of the form  $\sim n^{-\alpha/2} \sqrt{\ln n}$  at worst (this happens when  $A$  is a random set but we think it is a face).

To find a particular suitable  $p$  for a small set  $A$ , we note that if  $n^{-1} \ln |A| \leq \delta$  for some  $\delta < 1/4$ , the value of  $p$  obtained by substituting  $\rho = \sqrt{\delta}$  into the right-hand side of (4.5.2) would satisfy the inequality since the right-hand side is an increasing function of  $\rho$  and since  $\rho(A, p) \leq \sqrt{\delta}$  for any  $p$  by (4.5.1). It is interesting that an improvement in the cardinality estimate can be achieved by simply ignoring a (random) part of the information contained in the standard Hamming distance.

## 5 Corollaries, Remarks and Possible Ramifications

**5.1 Testing emerging exponential growth.** As an illustration, we show that our approach allows one to test in polynomial time whether the number of perfect matchings (Example 1.2) is exponentially large in the number of vertices of the graph (in the sense defined below). Let  $G = (V, E)$  be a graph with  $|V| = n$  vertices and let  $\mathcal{F}$  be the set of all perfect matchings in  $G$ . Economical embedding 2.4 allows us to identify

$\mathcal{F}$  with a subset  $F$  of the Boolean cube  $\{0, 1\}^m$  with  $m \leq n(\log_2 n + 1)$  and to construct efficiently Distance Oracle 2.2 for  $F$ . Suppose that  $|\mathcal{F}| \leq \exp\{n^\alpha\}$  for some  $\alpha < 1$ . Corollary 3.11 implies that  $\rho(F) = O(n^{\alpha/2}m^{-1/2})$ . Thus using Algorithm 3.3 to compute  $\Delta(F)$  within error  $\epsilon = \sqrt{n}$  and then Corollary 3.11 to interpret the result, we will be able to conclude that  $|\mathcal{F}| = O(\exp\{n^\beta\})$  for any  $\beta > (1 + \alpha)/2$ . Similarly, if  $|\mathcal{F}| \geq \exp\{n^\alpha\}$  for some  $\alpha > 0$ , our method would allow us to conclude for any  $\beta < 2\alpha - 1$  that  $|\mathcal{F}| = \Omega(\exp\{n^\beta\})$ . The estimate is, of course, void for  $\alpha < 1/2$  but it improves as  $\alpha$  approaches 1. For example, if  $|\mathcal{F}|$  has the order of  $\exp\{n^{0.95}\}$ , our method would allow us to conclude that  $|\mathcal{F}|$  is greater than  $\exp\{n^{0.89}\}$  and smaller than  $\exp\{n^{0.98}\}$ . Thus we can tell the order  $|\mathcal{F}| \approx \exp\{n^{0.99}\}$  from the order  $|\mathcal{F}| \approx \exp\{n\}$  and to distinguish them we have to solve the minimum weight matching problem a constant number of times. Bounds of this type for perfect matchings in general graphs are new.

Similarly, we can test whether the number of colored spanning subgraphs is exponentially large in the number  $r$  of colors (see Example 1.4).

Implementing our approach, Ryckman [R] wrote an experimental C++ code to estimate the number of perfect matchings in a given bipartite graph (or, equivalently, to estimate the permanent of a 0-1 matrix). In the case of a bipartite graph, Optimization Oracle 1.1 is especially easy to construct. In this case, the corresponding problem, known as the *Assignment Problem*, is not only “theoretically easy”, but in practice large instances are routinely solved as particular cases of the minimum cost network flow problem, see for example Section 11.2 of [PS]. Theoretically, the algorithm can not compete in precision with the recent polynomial time approximation scheme of [JSV]. However, in practice the code appears to be working extremely fast on fairly large graphs (it was tested on graphs with up to 256 vertices) and produces estimates which, although crude, are often non-trivial. Generally speaking, we think that our approach can be useful for problems of large size where some fairly crude estimates of the cardinality are needed and where the underlying optimization problem is especially easy (as in Example 1.3).

**5.2 Connections to Monte Carlo methods.** The main idea of our approach can be described as follows: given a (finite) ambient space  $\Omega$  and a set  $A \subset \Omega$ , we estimate the cardinality  $|A|$  by choosing a certain distance function  $d$  in  $\Omega$  and estimating the average distance

$$\Delta(A) = \frac{1}{|\Omega|} \sum_{x \in \Omega} d(x, A) \quad \text{where } d(x, A) = \min_{y \in A} d(x, y)$$

from  $x \in \Omega$  to  $A$ . We get the classical Monte Carlo method if the distance function  $d$  is the simplest possible:

$$d(x, y) = \begin{cases} 1 & \text{if } x \neq y \\ 0 & \text{if } x = y. \end{cases}$$

In this case,  $\Delta(A) = 1 - |A|/|\Omega|$ , so there is a direct relation between  $\Delta(A)$  and  $|A|$ . It is well understood that the main difficulty with the Monte Carlo method is that if  $|A|$  is “exponentially small” compared to  $|\Omega|$  then to get a non-trivial bound for  $|A|$ , we have to compute  $\Delta(A)$  with exponentially high precision. In this paper, we showed that in many interesting cases one can choose a different distance function  $d$ , so that the distance  $d(x, A)$  from a point  $x \in \Omega$  to  $A$  is efficiently computable and to get a meaningful estimate of  $|A|$  even for exponentially small sets  $A$ , one needs to compute  $\Delta(A)$  with a polynomial precision. Hence our approach can be considered as a natural extension of the Monte Carlo method.

In this paper, we choose  $\Omega$  to be the Boolean cube endowed either with the standard Hamming distance (section 3) or with its randomized version (section 4). In many cases, other embeddings might be of interest. In Example 1.3, all subsets  $Y \in \mathcal{F}$  have the same cardinality  $k$ . Straightforward embedding 2.3 identifies  $\mathcal{F}$  with a subset  $F$  of the section

$$\Omega = \left\{ (\xi_1, \dots, \xi_n) \in \{0, 1\}^n : \sum_{i=1}^n \xi_i = k \right\}$$

of the Boolean cube  $\{0, 1\}^n$ . Tighter estimates for  $|F|$  can be obtained by considering  $F$  to be a subset of  $\Omega$  and not of the whole cube  $\{0, 1\}^n$ .

In Example 1.4, the set  $\mathcal{F}$  of properly colored spanning subgraphs can be naturally identified with a subset  $F$  of the direct product  $\Omega = E_1 \times \dots \times E_r$ . Let  $d_i$  be a distance function on  $E_i$  for  $i = 1, \dots, r$  and let  $d$  be the corresponding  $L^1$  distance on  $\Omega$ . One can show that Distance Oracle 2.2 can be efficiently constructed for any choice of  $d_i$ . How should we choose  $d_i$  to get the best possible estimates for the cardinality  $|F|$ ? Note, that looking for such  $d_i$ , we are trying to satisfy two competing requirements: neighborhoods of “small” sets  $F \subset \Omega$  should be as small as possible while neighborhoods of “large” sets  $F \subset \Omega$  should be as large as possible. Perhaps one should use a whole family of distance functions  $d_i$  and combine the resulting estimates. The general isoperimetric inequality of [ABS] can be very useful for investigating that. Note, that economical embedding 2.4 results in choosing  $d_i$  so that  $E_i$  becomes isometric to a subset of the Boolean cube  $\{0, 1\}^{m_i}$  with  $m_i = \lceil \log_2 |E_i| \rceil$ .

**5.3 Weighted counting.** Let  $\mathcal{F} \subset 2^X$  be a family of subsets of the ground set  $X = \{1, \dots, n\}$  and let  $\mu(i) = p_i/q_i > 0$  be a rational weight of  $i \in X$ , where  $p_i, q_i \in \mathbb{N}$ . Let us define

$$\mu(Y) = \prod_{i \in Y} \mu(i) \text{ for } Y \in \mathcal{F} \quad \text{and} \quad \mu(\mathcal{F}) = \sum_{Y \in \mathcal{F}} \mu(Y).$$

We may be interested to estimate  $\mu(\mathcal{F})$ . There are several ways to extend our methods to problems of this type, here we sketch one. For every  $i \in X$ , let  $m_i = \lceil \log_2(p_i + q_i) \rceil$ . Let us choose subsets  $A_i \subset C_{m_i}$  and  $B_i \subset C_{m_i}$  such that  $|A_i| = p_i$ ,  $|B_i| = q_i$  and  $A_i \cap B_i = \emptyset$ . Let  $m = m_1 + \dots + m_n$  and let us identify

$$C_m = C_{m_1} \times \dots \times C_{m_n}.$$

For  $Y \subset \mathcal{F}$  let  $Z_Y \subset C_m$  be the direct product of  $n$  factors, the  $i$ -th factor being  $A_i$  if  $i \in Y$  and  $B_i$  if  $i \notin Y$ . Finally, let  $F \subset C_m$  be the union of all  $Z_Y$  for  $Y \in \mathcal{F}$ . We see that  $\mu(\mathcal{F}) = (q_1 \cdots q_n)^{-1} |F|$ . Moreover, one can define subsets  $A_i$  and  $B_i$  in such a way that Optimization Oracle 1.1 for  $\mathcal{F}$  gives rise to Distance Oracle 2.2 for  $F$ . This construction corresponds to the straightforward embedding 2.3. In some cases, there is a way to come up with an economical embedding in the spirit of 2.4.

### 6 Proofs of Theorems 4.4 and 4.5

**DEFINITION 6.1.** We recall that  $C_N$  is the Boolean cube  $\{0, 1\}^N$  endowed with the uniform probability measure and that  $\Lambda_N$  is the Boolean cube  $\{0, 1\}^N$  endowed with the probability measure of Definition 4.1. Let  $\Omega_N = C_N \times \Lambda_N$ . We consider the product measure on  $\Omega_N$ , so

$$\mathbf{P} \{(x, l)\} = p^{|l|} q^{N-|l|} 2^{-N} \quad \text{where } |l| = \lambda_1 + \dots + \lambda_N \text{ for } l = (\lambda_1, \dots, \lambda_N).$$

Hence a point  $(x, l) \in \Omega_N$  is interpreted as a vector of  $2N$  independent random variables  $(\xi_1, \dots, \xi_N; \lambda_1, \dots, \lambda_N)$ , where  $\mathbf{P} \{\xi_i=0\} = \mathbf{P} \{\xi_i=1\} = 1/2$ ,  $\mathbf{P} \{\lambda_i = 1\} = p$  and  $\mathbf{P} \{\lambda_i = 0\} = q$ . We observe that

$$\Delta(A, p) = \mathbf{E} d_l(x, A). \tag{6.1.1}$$

**LEMMA 6.2.** *Let  $A \subset C_N$  be a set. Then for every  $\delta \geq 0$*

$$\mathbf{P} \{(x, l) \in \Omega_N : |d_l(x, A) - \Delta(A, p)| \geq \delta\} \leq 2e^{-\delta^2/N}.$$

*Proof.* Since  $d_l(x, A)$  is a function of  $2N$  independent random variables, the proof follows by Lemma 3.4. □

Next, we need an analogue of the scaling trick 3.5.

**LEMMA 6.3.** *Let us fix positive integers  $k$  and  $n$  and let  $N = kn$ . Let us identify  $C_N = (C_n)^k$ ,  $\Lambda_N = (\Lambda_n)^k$  and  $\Omega_N = (\Omega_n)^k$ . Thus a point*

$(x, l) \in \Omega_N$  is identified with  $x = (x_1, \dots, x_k; l_1, \dots, l_k)$ , where  $x_i \in C_n$  and  $l_i \in \Lambda_n$ .

For a subset  $A \subset C_n$ , let  $B = A^k \subset C_N$ . Then

$$d_l(x, B) = \sum_{i=1}^k d_{l_i}(x_i, A) \quad \text{and} \quad \Delta(B, p) = k\Delta(A, p).$$

*Proof.* Clearly,

$$d_l(x, y) = \sum_{i=1}^k d_{l_i}(x_i, y_i) \quad \text{for all } x, y \in C_N$$

and the first identity follows. Now, by (6.1.1)

$$\Delta(B, p) = \mathbf{E} d_l(x, B) = \sum_{i=1}^k \mathbf{E} d_{l_i}(x_i, A) = k\Delta(A, p). \quad \square$$

Now we are ready to prove Theorem 4.4.

*Proof of Theorem 4.4.* Let  $N = nk$  and let us identify  $C_N = (C_n)^k$ ,  $\Lambda_N = (\Lambda_n)^k$  and  $\Omega_N = (\Omega_n)^k$ . Let  $B = A^k \subset C_N$  as in Lemma 6.3. Applying Lemma 6.2, we get

$$\mathbf{P} \left\{ (x, l) \in \Omega_N : |d_l(x, B) - \Delta(B, p)| \geq \delta \right\} \leq 2e^{-\delta^2/N}$$

for any  $\delta \geq 0$ . Using Lemma 6.3, we conclude:

$$\mathbf{P} \left\{ (x, l) \in \Omega_N : \left| \frac{1}{k} \sum_{i=1}^k d_{l_i}(x_i, A) - \Delta(A, p) \right| \geq \delta/k \right\} \leq 2e^{-\delta^2/N}.$$

Let us choose  $\delta = \epsilon k$ . Hence

$$\mathbf{P} \left\{ (x, l) \in \Omega_N : \left| \frac{1}{k} \sum_{i=1}^k d_{l_i}(x_i, A) - \Delta(A, p) \right| \geq \epsilon \right\} \leq 2e^{-\epsilon^2 k/n}.$$

Since  $k \geq 3n/\epsilon^2$ , the proof follows. □

Next, we need a (crude) version of inequality (3.7.1).

LEMMA 6.4. Let  $\epsilon \geq 0$ , let  $r(\epsilon) = pN(1 - \epsilon)/2$ . Let  $y \in C_N$  be a point. Then

$$\mathbf{P} \left\{ (x, l) \in \Omega_N : d_l(x, y) \leq r(\epsilon) \right\} \leq e^{-\epsilon^2 pN/4}.$$

*Proof.* Without loss of generality we may assume that  $y = 0$ . Then

$$\mathbf{P} \left\{ (x, l) \in \Omega_N : d_l(x, 0) \leq r(\epsilon) \right\} = \mathbf{P} \left\{ (x, l) \in \Omega_N : \sum_{i=1}^N \xi_i \lambda_i \leq r(\epsilon) \right\},$$

where  $x = (\xi_1, \dots, \xi_N)$  and  $l = (\lambda_1, \dots, \lambda_N)$ . Let  $\zeta_i = \xi_i \lambda_i$ . Then  $\zeta_i$ ,  $i = 1, \dots, N$  are independent random variables such that  $\mathbf{P} \{ \zeta_i = 1 \} = p/2$

and  $\mathbf{P} \{ \zeta_i = 0 \} = 1 - p/2$ . Hence

$\mathbf{P} \{ (x, l) \in \Omega_N : d_l(x, y) \leq r(\epsilon) \} = \mathbf{P} \{ \zeta_1 + \dots + \zeta_N \leq r(\epsilon) \} \leq e^{-\epsilon^2 p N/4}$   
 by a corollary of Hoeffding’s inequality (see Corollary 5.6 of [M]).  $\square$

Now we are ready to prove the first part of Theorem 4.5.

*Proof of inequality (4.5.1).* Let us choose a positive integer  $m$ , let  $N = mn$ , let  $C_N = (C_n)^m$ , and let  $\Lambda_N = (\Lambda_n)^m$ . Let  $B = A^m \subset C_N$  as in Lemma 6.3.

Let us choose an  $\alpha > 0$ . Applying Lemma 6.4, we obtain

$\mathbf{P} \{ (x, l) \in \Omega_N : d_l(x, B) \leq pN(1 - \sqrt{\alpha})/2 \} \leq |B|e^{-\alpha p N/4} = (|A|e^{-\alpha p n/4})^m$ .  
 Therefore, by Lemma 6.3

$$\mathbf{P} \left\{ (x, l) \in \Omega_N : \frac{1}{m} \sum_{i=1}^m d_{l_i}(x_i, A) \leq pn(1 - \sqrt{\alpha})/2 \right\} \leq (|A|e^{-\alpha p n/4})^m.$$

The right-hand side of the inequality tends to 0 provided  $\alpha > 4 \ln |A|/pn$ . Since by the Law of Large Numbers

$$\frac{1}{m} \sum_{i=1}^m d_{l_i}(x_i, A) \longrightarrow \Delta(A, p) \quad \text{in probability as } m \rightarrow +\infty,$$

we must have

$$\Delta(A, p) \geq pn(1 - \sqrt{\alpha})/2 \quad \text{for any } \alpha > 4 \ln |A|/pn.$$

Hence

$$\Delta(A, p) \geq pn(1 - \sqrt{\alpha})/2 \quad \text{for } \alpha = 4 \ln |A|/pn,$$

which is equivalent to (4.5.1).  $\square$

In section 3, we used the sharp isoperimetric inequality (Theorem 3.8) for the Hamming distance in  $C_n$  to get a sharp upper bound for  $n^{-1} \log_2 |A|$ . Unfortunately, we don’t know of a similar result for the randomized Hamming distance. To prove (4.5.2)–(4.5.3), we proceed by induction on  $n$  in a way resembling that of [T] (see also Remark 6.9).

We start with a simple technical result.

LEMMA 6.5. *For any  $0 \leq \epsilon \leq 1$ , any  $\gamma \geq 0$  and any  $0 < p \leq 1$  and  $q = 1 - p$  we have*

$$\begin{aligned} \min \left\{ \frac{p\gamma}{2} + \ln \frac{1}{1 + \epsilon}, p \ln \frac{1}{1 - \epsilon} + q \ln \frac{1}{1 + \epsilon} \right\} \\ \leq \max \left\{ 0, \ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \ln 2 \right\}. \end{aligned}$$

*Proof.* Fixing  $p, q$  and  $\gamma$ , let

$$f(\epsilon) = \frac{p\gamma}{2} + \ln \frac{1}{1 + \epsilon} \quad \text{and} \quad g(\epsilon) = p \ln \frac{1}{1 - \epsilon} + q \ln \frac{1}{1 + \epsilon}.$$



Then  $f(0) \geq 0$  and  $f(\epsilon)$  is decreasing whereas  $g(\epsilon)$  behaves as follows:  $g(0) = 0$  and if  $p \geq q$  then  $g(\epsilon)$  is increasing and if  $p < q$  then  $g(\epsilon)$  is decreasing for  $0 < \epsilon < q - p$  and increasing for  $q - p < \epsilon < 1$ . Furthermore,  $f(\epsilon_0) = g(\epsilon_0)$  at the single point  $\epsilon_0 = (e^{\gamma/2} - 1)/(1 + e^{\gamma/2})$ , where  $f(\epsilon_0) = g(\epsilon_0) = \ln(1 + e^{\gamma/2}) - q\gamma/2 - \ln 2$ . The proof now follows.  $\square$

DEFINITION 6.6. Let  $\mu_n$  (or simply  $\mu$ ) denote the uniform probability measure in  $C_n$ . Hence  $\mu(A) = |A|/2^n$ .

The induction is based on the following lemma.

LEMMA 6.7. Let  $A \subset C_{n+1}$  be a set. Let

$$A_0 = \{x \in C_n : (x, 0) \in A\} \quad \text{and} \quad A_1 = \{x \in C_n : (x, 1) \in A\}.$$

For  $l \in \Lambda_n$  let  $(l, 0) \in \Lambda_{n+1}$  denote  $l$  appended by  $\lambda_{n+1} = 0$  and let  $(l, 1) \in \Lambda_{n+1}$  denote  $l$  appended by  $\lambda_{n+1} = 1$ . Let

$$\Delta_0(A, p) = \mathbf{E} d_{(l,0)}(x, A) \quad \text{and} \quad \Delta_1(A, p) = \mathbf{E} d_{(l,1)}(x, A),$$

where the expectation is taken with respect to a random  $(x, l) \in C_{n+1} \times \Lambda_n$ .

Then

$$\frac{\mu_n(A_0) + \mu_n(A_1)}{2} = \mu_{n+1}(A); \tag{6.7.1}$$

$$\Delta(A, p) = q\Delta_0(A, p) + p\Delta_1(A, p); \tag{6.7.2}$$

$$\Delta_0(A, p) \leq \Delta(A_i, p) \quad \text{for } i = 0, 1; \tag{6.7.3}$$

$$\Delta_1(A, p) \leq \Delta(A_i, p) + \frac{1}{2} \quad \text{for } i = 0, 1; \tag{6.7.4}$$

$$\Delta_1(A, p) \leq \frac{\Delta(A_0, p) + \Delta(A_1, p)}{2}. \tag{6.7.5}$$

*Proof.* The proof is straightforward, cf. also proof of Lemma 2.1.2 of [T] and proof of Lemma 2.5 of [B].  $\square$

Now we use induction to get a preliminary bound.

LEMMA 6.8. Suppose that for some  $\gamma \geq 0$ ,  $0 < p \leq 1$  and  $q = 1 - p$ ,

$$\ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \ln 2 \geq 0.$$

Then for any non-empty set  $A \subset C_n$  we have

$$\gamma\Delta(A, p) + \ln \mu(A) \leq n \left( \ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \ln 2 \right).$$

*Proof.* We proceed by induction on  $n$ . If  $n = 1$  then two cases are possible:

$A$  consists of a single point,  $\mu(A) = 1/2$  and  $\Delta(A, p) = p/2$  (see Example 4.2);

$$A = \{0, 1\}, \mu(A) = 1 \text{ and } \Delta(A, p) = 0.$$

In either case, the inequality holds.

Suppose that the inequality holds for non-empty subsets of  $C_n$ . Let us prove that it holds for non-empty  $A \subset C_{n+1}$ . Let us define  $A_0, A_1 \subset C_n$  as in Lemma 6.7. From (6.7.1) it follows that either

$$\mu_n(A_0) = (1 - \epsilon)\mu_{n+1}(A) \quad \text{and} \quad \mu_n(A_1) = (1 + \epsilon)\mu_{n+1}(A)$$

or

$$\mu_n(A_1) = (1 - \epsilon)\mu_{n+1}(A) \quad \text{and} \quad \mu_n(A_0) = (1 + \epsilon)\mu_{n+1}(A)$$

for some  $0 \leq \epsilon \leq 1$ .

Let  $B$  be the one of the sets  $A_0, A_1$  that has a bigger measure  $\mu_n$  (either of the two if  $\mu_n(A_0) = \mu_n(A_1)$ ) and let  $D$  be the one of the sets  $A_0, A_1$  that has a bigger value of  $\Delta(\cdot, p)$  (either of the two if  $\Delta(A_0, p) = \Delta(A_1, p)$ ). Then

$$\mu_n(B) \geq (1 + \epsilon)\mu_{n+1}(A) \quad \text{and} \quad \mu_n(D) \geq (1 - \epsilon)\mu_{n+1}(A).$$

Furthermore, by (6.7.3)

$$\Delta_0(A, p) \leq \Delta(B, p) \quad \text{and} \quad \Delta_0(A, p) \leq \Delta(D, p)$$

whereas by (6.7.3) and (6.7.5)

$$\Delta_1(A, p) \leq \Delta(B, p) + \frac{1}{2} \quad \text{and} \quad \Delta_1(A, p) \leq \Delta(D, p).$$

Hence we get

$$\gamma\Delta_0(A, p) + \ln \mu_{n+1}(A) \leq \gamma\Delta(B, p) + \ln \mu_n(B) + \ln \frac{1}{1 + \epsilon}$$

and

$$\gamma\Delta_1(A, p) + \ln \mu_{n+1}(A) \leq \min \left\{ \gamma\Delta(B, p) + \ln \mu_n(B) + \ln \frac{1}{1 + \epsilon} + \frac{\gamma}{2}, \right. \\ \left. \gamma\Delta(D, p) + \ln \mu_n(D) + \ln \frac{1}{1 - \epsilon} \right\}.$$

Clearly,  $B$  is non-empty. Assume first, that  $D$  is non-empty as well. Applying the induction hypothesis to  $B$  and  $D$ , we conclude that

$$\gamma\Delta_0(A, p) + \ln \mu_{n+1}(A) \leq n \left( \ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \ln 2 \right) + \ln \frac{1}{1 + \epsilon}$$

and

$$\gamma\Delta_1(A, p) + \ln \mu_{n+1}(A) \\ \leq n \left( \ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \ln 2 \right) + \min \left\{ \ln \frac{1}{1 + \epsilon} + \frac{\gamma}{2}, \ln \frac{1}{1 - \epsilon} \right\}.$$

Adding the first inequality multiplied by  $q$  and the second inequality multiplied by  $p$  and using (6.7.2), we get

$$\begin{aligned} \gamma\Delta(A, p) + \ln \mu_{n+1}(A) &\leq n \left( \ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \ln 2 \right) \\ &\quad + \min \left\{ \frac{p\gamma}{2} + \ln \frac{1}{1 + \epsilon}, p \ln \frac{1}{1 - \epsilon} + q \ln \frac{1}{1 + \epsilon} \right\}. \end{aligned}$$

The desired inequality follows by Lemma 6.4.

If  $D$  is empty then  $\mu_n(B) = 2\mu_{n+1}(A)$  and we obtain

$$\gamma\Delta_0(A, p) + \ln \mu_{n+1}(A) \leq \gamma\Delta(B, p) + \ln \mu_n(B) - \ln 2$$

and

$$\gamma\Delta_1(A, p) + \ln \mu_{n+1}(A) \leq \gamma\Delta(B, p) + \ln \mu_n(B) - \ln 2 + \frac{\gamma}{2}$$

Adding the first inequality multiplied by  $q$  to the second inequality multiplied by  $p$  and using (6.7.2) and the induction hypothesis, we get

$$\begin{aligned} \gamma\Delta(A, p) + \ln \mu_{n+1}(A) &\leq \gamma\Delta(B, p) + \ln \mu_n(B) - \ln 2 + \frac{p\gamma}{2} \\ &\leq n \left( \ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \ln 2 \right) + \left( \frac{\gamma}{2} - \frac{q\gamma}{2} - \ln 2 \right) \\ &\leq (n + 1) \left( \ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \ln 2 \right), \end{aligned}$$

which completes the proof. □

Now we are ready to complete the proof of Theorem 4.5.

*Proof of (4.5.2)–(4.5.3).* By Lemma 6.8,

$$\frac{\ln |A|}{n} = \frac{\ln \mu_n(A)}{n} + \ln 2 \leq \ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \frac{\gamma\Delta(A, p)}{n} = \ln(1 + e^{\gamma/2}) - \frac{\gamma}{2} + \gamma\rho$$

provided

$$\ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \ln 2 \geq 0.$$

We optimize the inequality on  $\gamma \geq 0$ . Let

$$\gamma = 2 \ln \left( \frac{1}{2\rho} - 1 \right).$$

Since we assumed that  $\rho \leq 1/4$ , we have  $\gamma \geq 0$ . Furthermore,

$$\begin{aligned} \ln(1 + e^{\gamma/2}) - \frac{q\gamma}{2} - \ln 2 &= \ln \frac{1}{2\rho} - q \ln \left( \frac{1}{2\rho} - 1 \right) - \ln 2 \\ &= -\ln(1 - 2\rho) + p(\ln(1 - 2\rho) - \ln(2\rho)) - \ln 2 \geq 0, \end{aligned}$$

because of (4.5.2). Therefore,

$$\frac{\ln |A|}{n} \leq \ln \frac{1}{2\rho} - \ln \frac{1 - 2\rho}{2\rho} + 2\rho \ln \frac{1 - 2\rho}{2\rho} = 2\rho \ln \frac{1}{2\rho} + (1 - 2\rho) \ln \frac{1}{1 - 2\rho}$$

and (4.5.3) follows. □

**REMARK 6.9.** Our proof of (4.5.2)–(4.5.3) can be considered as an “additive” version of Talagrand’s method [T]. Indeed, Talagrand’s approach

very roughly can be stated as follows. Let  $\Omega$  be a space with the distance function  $d$  and probability measure  $\mu$ . To prove an isoperimetric inequality for  $A \subset \Omega$ , we first find a uniform bound for the expression  $\mu^\alpha(A) \cdot \mathbf{E} \exp\{\tau d(x, A)\}$  and then adjust parameters  $\alpha > 0$  and  $\tau > 0$ . This way tight inequalities are obtained in [T] for sets  $A$  of large measure, most often with  $\mu(A) \geq 1/2$ . We are mostly interested in sets of a small measure. One can check that for “small sets”  $A$  the inequalities of [T] are very far from sharp, which is, of course, should not be perceived as a “fault” of the method, since the method was designed for totally different problems. We find a uniform bound for the expression  $\ln \mu(A) + \tau \mathbf{E} d(x, A)$ , which looks like Talagrand’s functional with “exp” removed. Our method seems to produce reasonably good bounds for small sets  $A$  but it fails miserably for large  $A$ , with  $\mu(A) = 1/2$ , say. As should have been expected, the case of “middle-sized” sets is the most complicated.

**Acknowledgment.** The authors are grateful to M. Gromov, G. Kalai and B. Sudakov for many helpful discussions and to the anonymous referee for suggestions.

### References

- [ABS] N. ALON, R. BOPPANA, J. SPENCER, An asymptotic isoperimetric inequality, *Geom. Funct. Anal.* 8 (1998), 411–436.
- [B] A. BARVINOK, Approximate counting via random optimization, *Random Structures & Algorithms* 11:2 (1997), 187–198.
- [GLS] M. GRÖTSCHEL, L. LOVÁSZ, A. SCHRIJVER, *Geometric Algorithms and Combinatorial Optimization*, second ed., *Algorithms and Combinatorics*, 2, Springer-Verlag, Berlin, 1993.
- [J] M. JERRUM, The computational complexity of counting, in “Proceedings of the International Congress of Mathematicians, Vol. 1,2 (Zürich, 1994)”, Birkhäuser, Basel (1995), 1407–1416.
- [JS1] M. JERRUM, A. SINCLAIR, Approximating the permanent, *SIAM J. Comput.* 18:6 (1989), 1149–1178.
- [JS2] M. JERRUM, A. SINCLAIR, The Markov chain Monte Carlo method: an approach to approximate counting and integration, in “Approximation Algorithms for NP-hard Problems” D.S. Hochbaum, ed., PWS, Boston (1997), 483–520.
- [JSV] M. JERRUM, A. SINCLAIR, E. VIGODA, A polynomial-time approximation algorithm for the permanent of a matrix with non-negative entries, *Electronic Colloquium on Computational Complexity*, Report TR00-079, <ftp://ftp.eccc.uni-trier.de/pub/eccc/reports/2000/TR00-079/index.html>

- [JVV] M. JERRUM, L.G. VALIANT, V.V. VAZIRANI, Random generation of combinatorial structures from a uniform distribution, *Theoret. Comput. Sci.*, 43:2-3 (1986), 169–188.
- [L] I. LEADER, Discrete isoperimetric inequalities, in “Probabilistic Combinatorics and its Applications (San Francisco, CA, 1991)”, Proc. Sympos. Appl. Math. 44, Amer. Math. Soc., Providence, RI (1991), 57–80.
- [Li] J.H. VAN LINT, Introduction to Coding Theory, Third edition, Graduate Texts in Mathematics 86, Springer-Verlag, Berlin, 1999.
- [LoP] L. LOVÁSZ, M.D. PLUMMER, Matching Theory, North-Holland Mathematics Studies 121, Annals of Discrete Mathematics 29, North-Holland Publishing Co., Amsterdam-New York; Akadémiai Kiadó (Publishing House of the Hungarian Academy of Sciences), Budapest, 1986.
- [M] C. MCDIARMID, On the method of bounded differences, in “Surveys in Combinatorics 1989 (Norwich, 1989)”, London Math. Soc. Lecture Note Ser. 141, Cambridge Univ. Press, Cambridge (1989), 148–188.
- [MiS] V.D. MILMAN, G. SCHECHTMAN, Asymptotic Theory of Finite-Dimensional Normed Spaces (with an Appendix by M. Gromov), Springer Lecture Notes in Mathematics 1200, 1986.
- [P] C.H. PAPADIMITRIOU, Computational Complexity, Addison-Wesley, Reading, Mass., 1994.
- [PS] C.H. PAPADIMITRIOU, K. STEIGLITZ, Combinatorial Optimization: Algorithms and Complexity, Dover, NY, 1998.
- [R] E.M. RYCKMAN, Code for permanent approximations, experimental C++ code, available at <http://www.math.lsa.umich.edu/~barvinok/papers.html>
- [T] M. TALAGRAND, Concentration of measure and isoperimetric inequalities in product spaces, *Inst. Hautes Études Sci. Publ. Math.* 81 (1995), 73–205.

ALEXANDER BARVINOK, Department of Mathematics, University of Michigan,  
Ann Arbor, MI 48109-1109, USA [barvinok@math.lsa.umich.edu](mailto:barvinok@math.lsa.umich.edu)

ALEX SAMORODNITSKY, Institute for Advanced Study, Einstein Drive, Princeton,  
NJ 08540, USA [asamor@ias.edu](mailto:asamor@ias.edu)

Submitted: June 2000

Revised version: January 2001