

Optimal Reconstruction of Graphs Under the Additive Model

Vladimir Grebinski, Gregory Kucherov

► **To cite this version:**

Vladimir Grebinski, Gregory Kucherov. Optimal Reconstruction of Graphs Under the Additive Model. [Research Report] RR-3171, INRIA. 1997. inria-00073517

HAL Id: inria-00073517

<https://hal.inria.fr/inria-00073517>

Submitted on 24 May 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

*Optimal Reconstruction of Graphs Under the
Additive Model*

Vladimir Grebinski and Gregory Kucherov

N° 3171

———— THÈME 2 ————



*R*apport
de recherche

Optimal Reconstruction of Graphs Under the Additive Model

Vladimir Grebinski and Gregory Kucherov

Thème 2 — Génie logiciel
et calcul symbolique
Projet Euréca

Rapport de recherche n3171 — Mai 1997 — 18 pages

Abstract: We study the problem of combinatorial search for graphs under the additive model. The main result concerns the reconstruction of *bounded degree* graphs, i.e. graphs with the degree of all vertices bounded by a constant d . We show that such graphs can be reconstructed in $O(dn)$ non-adaptive queries, that matches the information-theoretic lower bound. The proof is based on the technique of separating matrices. In particular, a new upper bound is obtained for d -separating matrices, that settles an open question stated by Lindström in [16]. Finally, we consider several particular classes of graphs. We show how an optimal non-adaptive solution of $O(n^2/\log n)$ queries for general graphs can be obtained.

Key-words: combinatorial search, non-adaptive algorithms, graph reconstruction

(Résumé : *tsvp*)

An extended abstract of this paper is to appear in the Proceedings of the 5th Israeli Symposium on Theory of Computing and Systems. IEEE Press, June 1997.

Reconstruction optimale de graphes dans le modèle additif

Résumé : Nous étudions le problème de recherche combinatoire de graphes dans le modèle additif. Le résultat principal porte sur la reconstruction de graphes à degré borné, c'est à dire ceux dont le degré de tous les noeuds ne dépassent pas une constante d . Nous démontrons qu'un tel graphe peut être reconstruit en $O(dn)$ requêtes non-adaptatives, ce qui correspond à la borne inférieure par la théorie d'information. La preuve est basée sur la technique de matrices de séparation pour lesquelles une nouvelle borne supérieure est obtenue qui résout la question ouverte posée par Lindström dans [16]. Finalement, nous considérons quelques classes particulières de graphes. Nous démontrons comment une solution optimale non-adaptative de $O(n^2/\log n)$ requêtes peut être obtenue pour les graphes généraux.

Mots-clé : recherche combinatoire, algorithmes non-adaptatifs

1 Introduction and Definitions

Combinatorial Search studies problems of the following general type: determine an unknown object by means of indirect questions about this object. Perhaps the most common example of combinatorial search is the variety of problems of determining one or several counterfeit coins in a set using scales of some kind. Many of these problems still lack an optimal general solution.

Each instance of a Combinatorial Search problem has two main components: a finite *domain of objects* \mathcal{M} and a *class of queries* \mathcal{Q} , which is a family of functions from the domain of objects to a domain \mathcal{A} of *answers*. Given \mathcal{M} and \mathcal{Q} , the combinatorial search problem is to find a sequence of queries (q_1, q_2, \dots, q_k) , $q_i \in \mathcal{Q}$, such that the sequence of answers $(q_1(x), q_2(x), \dots, q_k(x))$ uniquely identifies the object $x \in \mathcal{M}$. A method for choosing queries (q_1, q_2, \dots, q_k) is called a (*combinatorial*) *search algorithm*. The complexity measure of a search algorithm is the maximal number k of required queries over all $x \in \mathcal{M}$. This implies that we are concerned with *query complexity* only. Precise complexity bounds to combinatorial search problems can be rarely obtained. Instead, one is usually interested in the asymptotic complexity, when $|\mathcal{M}|$ tends to infinity.

Monographs [6, 2] present detailed accounts of numerous results on Combinatorial Search problems. Variants of these problems abound in different application domains. For example, paper [10] deals with a problem motivated by genome analysis. Note that Combinatorial Search is closely related to Learning Theory, where the general framework is similar, except possibly that there is usually an infinity of objects and one is looking not necessarily for the object itself but for its approximation according to a given distance function.

In general, the choice of q_i in the sequence (q_1, q_2, \dots, q_n) depends on the answers $(q_1(x), \dots, q_{i-1}(x))$ obtained “so far”. If this dependence exists, the algorithm is called *adaptive* (or sequential). Otherwise, when all the queries can be given before any answer is known, the algorithm is called non-adaptive (or predetermined). In this paper, we deal with non-adaptive algorithms. Although they are obviously less powerful in general, non-adaptive algorithms usually admit “nicer” mathematical formulations that allow to use more powerful mathematical methods. Besides, in many cases (including those considered in this paper) non-adaptive algorithms achieve the power of adaptiveness, that is reach the lower bound. Note also that in non-adaptive algorithms all queries can be made in parallel, which is useful in many applications.

For the non-adaptive case, we reformulate the combinatorial search problem as follows: find a minimal number of queries $q_1, q_2, \dots, q_n \in \mathcal{Q}$ such that for every $x, y \in \mathcal{M}$, there is q_i , $1 \leq i \leq n$, such that $q_i(x) \neq q_i(y)$.

In contrast to the coin weighing problem where objects of \mathcal{M} are just elements or subsets of elements of a given set, the objects may be of a more complex nature, such as graphs or partially ordered sets (see [1, 2]). In case of graphs, different combinatorial search problems can be raised. One may look for an unknown edge in a *given* graph by asking, for a subset of vertices, whether one of the edge’s endpoints (or both) belongs to the subset. A more general problem, considered in this paper, consists of determining an unknown graph of a given class. Here again, subsets of vertices are queried, but the answer returned characterizes

some property of the subgraph induced by the subset. Finally, the third type of problem is to check whether an unknown graph belongs to a given class without actually determining the graph. This problem, known as *property testing*, received much attention in connection with the study of *evasiveness* property (see [17]). Another approach to property testing, in the framework of probabilistic algorithms and approximation, was recently introduced in [8, 9].

It is clear that for the same object domain \mathcal{M} , different classes of queries \mathcal{Q} lead to combinatorial search problems of different type and different complexity. Under the *additive model*, the domain of answers is the ring of integers \mathbb{Z} . This model is also called *quantitative*, as the queries \mathcal{Q} are usually about some quantitative property of the object. A typical example is to identify the subset of counterfeit coins using a spring scale under the knowledge of the difference in weight between a counterfeit and authentic coin (which allows to determine the number of counterfeit coins in a subset by weighing this subset). We will come back to this example in Sect. 2. Some additive models of combinatorial search are studied in [12, 7, 11].

In this paper we consider the problem of *searching for a graph under the additive model* defined as follows. The domain of objects, denoted \mathcal{G}_n is a class of simple graphs with n vertices labelled by natural numbers $1, 2, \dots, n$. (A graph is simple if it does not contain loops and multiple edges.) The queries that we are allowed to make about $G \in \mathcal{G}_n$ are of the following form: For a subset $V \subseteq \{1, \dots, n\}$ of vertices, how many edges are there in G between vertices of V ? More formally, how many edges occur in the intersection $G \cap K_V$, where K_V is the complete graph with the set of vertices V ?

In this paper we develop new techniques of non-adaptive additive search for a graph. Our main result concerns the search for *bounded degree* graphs, i.e. graphs with the degree of all vertices bounded by a constant d . We prove that such a graph can be reconstructed within $O(dn)$ non-adaptive queries, which matches the information-theoretic lower bound. The key intermediate result shows that a bipartite graph can be reconstructed in $O(dn)$ non-adaptive queries provided that the degree of vertices *on one side* is bounded by d while no restriction on the other part is made (we call such graphs one-sided bounded degree bipartite graphs). We also show how an optimal non-adaptive solution of $O(n^2/\log n)$ queries for general graphs can be obtained.

The results show the power of the considered model, gained by the possibility of testing a *set* of vertices and *counting* the number of edges between them. For comparison, if we are allowed to query only two vertices (that is, test one edge at a time), $\Omega(n^2)$ queries are needed for many natural classes of graphs, such as trees, matchings, hamiltonian cycles and paths, and others (see [2]).

The paper is organized as follows. Sect. 2 is devoted to separating matrices – the main tool for constructing non-adaptive additive algorithms. We consider two classes of separating matrices – so called d -separating and d -detecting matrices. Our main contribution here is a probabilistic proof of a new upper bound for d -separating matrices which matches modulo a constant factor the information-theoretic lower bound. In Sect. 3 we turn to our main subject of interest – searching for graphs under the additive model. We consider bounded-

degree graphs and prove that such graphs can be reconstructed in $O(dn)$ queries. Finally, in Sect. 4 we consider several particular classes of graphs. For some of them, that are not subclasses of degree bounded graphs, we show that our technique still applies. For others, we show that the constant factor can be improved.

2 Separating Matrices

Consider the following setting. Assume we have a set of *items* and each of them is assigned an integer value. Assume that we want to reconstruct the values by making queries about subsets of items. As noted in the introduction, this type of search problems is very common and is called *combinatorial group testing*. Note that each query can be associated to a $(0, 1)$ -vector q which is the incidence vector of the corresponding subset. Assume further that the result of a query is the sum of item values of the corresponding subset. This assumption typically corresponds to the *additive model* discussed in the introduction. It implies that if v is the vector of item values, the query result is the scalar product $\langle q, v \rangle$. Let us now restrict ourselves to non-adaptive algorithms. Then the whole algorithm can be represented by a $(0, 1)$ -matrix where each row is a query vector and each column corresponds to an item. This leads us to the following notion.

Definition 2.1 A $k \times n$ $(0, 1)$ -matrix M is called separating for a finite set of integer vectors $V \subseteq \mathbb{Z}^n$ iff for every $v_1, v_2 \in V$, $Mv_1 \neq Mv_2$ provided that $v_1 \neq v_2$.

In this section we study two important subclasses of separating matrices.

2.1 Optimal d -Separating Matrices

Definition 2.2 For a constant $d \in \mathbb{N}$, a d -separating matrix is a separating matrix for the set of $(0, 1)$ -vectors containing at most d entries equal to 1.

Equivalently, for a d -separating matrix, all the sums of any up to d columns are distinct. If a matrix has n columns, there will be $\sum_{i=0}^d \binom{n}{i}$ different sums of at most d columns. Since each entry of such a sum is at most d , a d -separating matrix has at least $\log_{d+1} \binom{n}{d} = (1 + o(1))d \log_d(n)$ rows (recall that d is a constant). A better lower bound was proved by Noga Alon [3]. Using the second moment method (cf. [4]), it was shown in [3] that there exists an absolute constant c such that for every $n > d$ any d -separating matrix has at least $\frac{2d}{\log d + c} \log(n/d)$ rows.

By definition, a matrix is 1-separating iff all its columns are different. Clearly, a 1-separating matrix with n columns and $\lceil \log_2 n \rceil$ rows can be easily constructed by setting the columns to be the binary representations of numbers $1, 2, \dots, n$. This matrix corresponds to the *non-adaptive binary search*, as it provides a non-adaptive analogue to the binary search procedure. For an arbitrary constant d , it is known that a d -separating matrix with asymptotically $d \log_2 n$ rows can be constructed (see [16] and [2, exercise 2.3.5]). For $d = 2$, the lower bound $(5/3) \log_2 n$ has been proved by Lindström [15], while no better upper bound

than $2 \log_2 n$ is known. This suggests that settling the multiplicative factor for the case of arbitrary constant d is difficult.

In this section we give a probabilistic proof that there exists a d -separating matrix which asymptotically meets the lower bound $\Omega(d \log_d n)$ up to a multiplicative constant independent of d . More precisely, we obtain the upper bound which is within the factor two of the lower bound $(2 + o(1))d \log_d n$ from [3]. This answers the question whether the upper bound $d \log_2 n$ can be improved, posed by Lindström in [16].

Before proceeding to the proof, we note a straightforward connection between d -separating matrices and a classical problem of *counterfeit coins*. A d -separating matrix with n columns solves the following problem. *Suppose we have n coins of which at most d are counterfeit. We are allowed to ask how many counterfeit coins occur in a subset. Find an optimal non-adaptive algorithm that determines all counterfeit coins.*

We now prove the main result of this section.

Theorem 1 (d -separating matrix) *For fixed d , there exists a d -separating matrix with n columns and asymptotically $4d \log_d n$ rows.*

Proof: We show that a random matrix M with asymptotically $4d \log_d n$ rows is d -separating with a positive probability. This will imply that such a matrix exists [4].

Consider a $(0, 1)$ -matrix M . Let A and B be two different subsets of columns, each of size at most d . We say that (A, B) is a conflicting pair for M iff the sum of columns A in M is equal to the sum of columns B . Without loss of generality, we assume (A, B) to have two additional properties:

1. $A \cap B = \emptyset$. Indeed, if (A, B) is a conflicting pair, then $(A \setminus B, B \setminus A)$ is a conflicting pair too.
2. A and B have the same size $0 < |A| = |B| \leq d$. This can be insured by adding to M an additional row with all entries equal to 1. This row will not be subject to the random choice of matrix' entries. Obviously, adding one row does not affect the asymptotic bound.

We are going to estimate the expected number of conflicting pairs in a random $k \times n$ matrix (k will be chosen later). We define

$$\chi(A, B, M) = \begin{cases} 1 & \text{if } (A, B) \text{ is a conflicting pair for } M \\ 0 & \text{otherwise} \end{cases}$$

We assume the uniform probabilistic distribution over the $n \times k$ $(0, 1)$ -matrices. This implies that each entry of M is chosen independently to be 0 or 1 with probability $1/2$. For each fixed i , $1 \leq i \leq d$, we consider all possible partitions σ of the set of $2i$ columns into two equal parts of i columns. There are clearly $\frac{1}{2} \binom{2i}{i}$ different partitions. For each partition σ , we consider all possible $\binom{n}{2i}$ subsets of $2i$ columns. Each subset R of $2i$ columns is thought to be split according to the partition σ above into two non-intersecting sets R_σ^+ and R_σ^- of size i . Compute now the expectation $\mathbb{E}(\chi(R_\sigma^+, R_\sigma^-, M))$ for fixed non-intersecting

R_σ^+ and R_σ^- , both of cardinality $i \leq d$. Clearly, the events “two sums of entries in a row are equal” are independent for different rows. Furthermore, since $A \cap B = \emptyset$, the probability for one row is $\sum_{j=0}^i \binom{i}{j} 2^{-i} = \binom{2i}{i} 2^{-2i} = \frac{1}{\sqrt{\pi i}} + O(\frac{1}{i})$. However, we will need an upper bound for all i , and we use that $\binom{2i}{i} 2^{-2i} \leq \frac{2}{3\sqrt{i}}$ for all $i \geq 1$. A way to show this is to prove by induction over $i \geq 1$ the stronger inequality $\binom{2i}{i} 2^{-2i} \leq \frac{2}{3} \frac{1}{e^{\frac{1}{8i}} \sqrt{i}}$ (consider the ratio of consecutive elements and use the fact that $e^{\frac{1}{8i}} \geq 1 + \frac{1}{8i}$). Using this, we have $\mathbb{E}(\chi(A, B, M)) = (\binom{2i}{i} 2^{-2i})^k \leq (\frac{2}{3})^k (\frac{1}{i})^{k/2}$.

Now the expected number of conflicting pairs can then be estimated as

$$\begin{aligned} \mathbb{E} \left(\sum_{i=1}^d \sum_{\sigma} \sum_R \chi(R_\sigma^+, R_\sigma^-, M) \right) &= \sum_{i=1}^d \binom{2i}{i} \binom{n}{2i} \left(\binom{2i}{i} 2^{-2i} \right)^k \leq \\ &\leq \sum_{i=1}^d 2^{2i} \left(\frac{en}{2i} \right)^{2i} \left(\frac{2}{3} \right)^k \left(\frac{1}{i} \right)^{k/2} = \sum_{i=1}^d \left(\frac{2}{3} \right)^k \frac{(en)^{2i}}{i^{2i+k/2}} \end{aligned}$$

An elementary analysis shows that when $k = 4d \log_d n$, every summand is $o(1)$ and the whole sum can be made smaller than 1. Since the expected value of the number of conflicting pairs is smaller than 1, there exists at least one matrix M without conflicting pairs and therefore d -separating. \square

2.2 d -Detecting Matrices

In this section we consider another important class of separating matrices.

Definition 2.3 *Let d be a constant. A $k \times n$ $(0, 1)$ -matrix, with n columns, is called d -detecting iff it is separating for the set of n -vectors $\{0, \dots, d-1\}^n$.*

Let v_1, v_2, \dots, v_n be the columns of a $(0, 1)$ -matrix. Then this matrix is d -detecting iff all the sums $\sum_{i=1}^n \epsilon_i v_i$ ($\epsilon_i = 0, 1, \dots, d-1$) are different. Such a set of vectors is called *detecting* in [14], hence our terminology.

Given n and d , we are interested in d -detecting matrices with minimal number of rows. Let k be the number of rows. An information-theoretic reasoning gives the inequality $d^n \leq (dn)^k$, and the lower bound $\Omega(n/(1 + \log_d n))$ for k . The problem has been studied by several authors, and particularly by Bernt Lindström in a series of papers [13, 14, 16]. In [14] Lindström presents a construction of a detecting matrix, using the theory of Möbius functions. This construction gives a solution of order $2n/\log_d n$, although this was not explicitly pointed out by the author. Moreover, this bound turns out to be optimal, that will be stated in Lemma 1 below. We refer to Appendix A for a summary of main Lindström’s results from [13, 14], and their application to the construction of an optimal d -detecting matrix. Here we summarize this construction in the following theorem.

Theorem 2 *For fixed d , a d -detecting matrix can be effectively constructed with n columns and asymptotically $2n/\log_d n$ rows.*

In [16] Lindström concentrates on the case $d = 2$ for which he proposes a construction of a detecting matrix based on elementary methods (the construction is also described in [2, 6]). The matrix has asymptotically $2n/\log_2 n$ rows. Lindström also proves that the construction is optimal, that is the bound $2n/\log_2 n$ is the asymptotic lower bound. Further references for the case $d = 2$ can be found in [16].

It is interesting that the construction of Theorem 2 matches the asymptotic lower bound for d -detecting matrices.

Lemma 1 *For a fixed d , any d -detecting matrix with n columns has at least $2n/\log_d n$ rows asymptotically.*

The proof is given in Appendix B. It is a generalization of the proof for $d = 2$ [16] based on the method attributed to L.Mozer.

Similarly to d -separating matrices, d -detecting matrices have a natural interpretation in terms of “generalized counterfeit coins problem”. Assume we have n coins and an unknown arbitrary number of them are false. Assume further that we know the weight α of an authentic coin, and that the weight of each false coin takes one of the values $\alpha + \delta i$ for $i = 1, \dots, d - 1$. One can think of i (the overweight of a coin) as the “measure of falsity”. We are allowed to weigh subsets of coins and thus measure the overall overweight of a subset. Determine the false coins and their falsity by possibly minimal number of weighing. It is easily seen that finding a non-adaptive solution of the generalized counterfeit coins problem is directly translated to constructing a d -detecting matrix with minimal number of rows. Note that for $d = 2$ we get the counterfeit coins problem described in Sect. 2.1 but with an arbitrary non-fixed number of false coins.

3 Reconstructing Bounded Degree Graphs

Now we turn to our main subject of interest – the problem of graph reconstruction under the additive model. Let \mathcal{G}_n be a class of labelled undirected graphs with n vertices labelled by $\{1, 2, \dots, n\}$. We consider simple graphs, that is graphs without loops or multiple edges.

We address the following problem. Reconstruct an unknown graph $G \in \mathcal{G}_n$ by means of queries of the following type: For a subset $V \subseteq \{1, \dots, n\}$ of vertices, how many edges are there in the intersection $G \cap K_V$, where K_V is the complete graph with the set of vertices V ? Clearly, each query is simply associated with a subset of $\{1, \dots, n\}$.

Using the results on separating matrices presented in Sect. 2, we solve this problem for an important subclass of graphs, namely the *bounded degree graphs*. These are graphs with the degree of vertices bounded by some constant d . Bounded degree graphs are quite common objects and cover such classes as matchings, circles and paths, trees with bounded branching degree, etc. Property testing for bounded degree graphs was considered in [9]. In

this section we propose an asymptotically optimal (modulo a constant factor) predetermined search algorithm for this class.

We first prove an auxiliary result. Using the results of Sect. 2.1 and 2.2, we prove the existence of an optimal predetermined algorithm for the class of *one-sided bounded degree bipartite graphs*. Consider a bipartite graph $G = (V, W, E)$, where $V \cup W$ is the set of vertices, $V \cap W = \emptyset$, and $E \subseteq V \times W$. For a constant d , G is called a one-sided (d -)bounded degree graph if $\deg(v) \leq d$ for every vertex $v \in V$.

Assume that $|V| = |W| = n$. By assuming that every node of V has degree d , it is easy to estimate the number of such graphs from below as $\binom{n}{d}^n = \Omega((n/d)^{dn})$. Since the answer to a query has $nd + 1$ potential values, any search algorithm requires at least $\log_{dn+1}(n/d)^{dn} = dn(1 + o(1))$ queries. We now prove that this lower bound can be met, modulo a constant factor, by a non-adaptive algorithm.

Theorem 3 *For a constant d , there exists a non-adaptive search algorithm for the class of one-sided d -bounded degree bipartite graphs with n vertices on each side, that requires $8dn$ queries asymptotically.*

Proof: Consider a bipartite graph $G = (V, W, E)$, where $V \cup W$ is the set of vertices, $V \cap W = \emptyset$, $|V| = |W| = n$, and $E \subseteq V \times W$ is the set of edges. Assume that $\deg(v) \leq d$ for all $v \in V$. By definition, each query is associated with a couple (V', W') , $V' \subset V$, $W' \subset W$, and has the form: how many edges of G are between vertices of V' and W' ? (what is $|E \cap (V' \times W')|$?)

For a vertex $v \in V$ and a subset $W' \subseteq W$, denote by $\deg_{W'}(v)$ the number of vertices of W' adjacent to v . Note that $\deg_{W'}(v) \leq d$ and can be determined by one query.

Fix a vertex $v \in V$. According to Theorem 1, we can find all its adjacent vertices in W using $4d \frac{\log n}{\log d}$ queries (think of adjacent vertices as being “counterfeit”, all other vertices in W being “authentic”). Each query asks about $\deg_{W'}(v)$ for some subset $W' \subseteq W$. Let $W_1, W_2, \dots, W_k \subseteq W$ (k asymptotically to $4d \frac{\log n}{\log d}$) be these subsets. Since the algorithm is predetermined, the subsets W_1, W_2, \dots, W_k are independent of v . In other words, if for some $v \in V$ we know $\deg_{W_i}(v)$ for every W_i , we can reconstruct all the adjacent vertices of v in W .

Now fix some W_i . By Theorem 2, we can find a sequence of subsets V_1, \dots, V_l , with l asymptotically to $2n \frac{\log d}{\log n}$, such that the queries $\langle (W_i, V_j), j = 1, \dots, l \rangle$ allow to reconstruct $\deg_{W_i}(v)$ for every $v \in V$ (think of $\deg_{W_i}(v)$ as “the degree of falsity” of v).

Repeating this algorithm for each W_i , we can determine $\deg_{W_i}(v)$ for all $v \in V$ and all W_i via $4d \frac{\log n}{\log d} \cdot 2n \frac{\log d}{\log n} = 8dn$ queries asymptotically. This allows us to reconstruct the adjacent vertices in W of each $v \in V$, that is to reconstruct the whole graph. \square

The key argument of the proof is that the algorithm implied by Theorem 1 is non-adaptive, i.e. the sets W_i don't depend on vertices of V . The fact that the algorithm implied by Theorem 2 is also non-adaptive does not affect the complexity bound but insures that the resulting algorithm is completely predetermined too. Specifically, it insures that

all the sets V_j are predetermined, and therefore all the queries (W_i, V_j) are mutually independent and can be made in any order.

We now use Theorem 3 to construct a separating set of queries for general bounded degree graphs. We start with computing the information-theoretic lower bound and for that we estimate from below the number of graphs with bounded degree. Instead of counting all such graphs, we will count a subclass of them, and show that their number is already sufficiently big.

Denote by $D(n, d)$ the set of labelled bipartite graphs with n vertices on each side with the degree of each vertex equal to d (d constant). This d -regular graph is a union of d disjoint matchings. Clearly, $|D(n, 0)| = 1$ and $|D(n, 1)| = n!$.

Consider a graph $G \in D(n, d)$. From G we can obtain a graph in $D(n, d + 1)$ by adding a matching which doesn't intersect with G . To estimate the number of possible extensions, consider the complement bipartite graph \bar{G} (an edge connecting the sides belongs to \bar{G} iff it does not belong to G). It is an $(n - d)$ -regular graph. Since the number of matching in a bipartite graphs is equal to the permanent of the adjacency matrix, from the Van der Waerden conjecture, proved by Egorychev and Falikman (see [5]), it follows that this graph has at least $(n - d)^n \frac{n!}{n^n}$ matchings. Obviously, none of them intersects with G .

On the other hand, consider a graph $G' \in D(n, d + 1)$. The number of matchings it contains is bounded from above by $(d + 1)^n$ (a better estimation is not important for our purposes). Thus, $|D(n, 0)| = 1$ and $|D(n, d)|(n - d)^n \frac{n!}{n^n} \leq |D(n, d + 1)|(d + 1)^n$. From this recurrence, $|D(n, d)| \geq \binom{n}{d}^n \left(\frac{n!}{n^n}\right)^d \geq (n/d)^{dn} (e^{-nd}) = \left(\frac{n}{ed}\right)^{nd}$. We obtain $\log_{n, d+1} |D(n, d)| \geq nd(1 + o(1))$, and thus any search algorithm for $D(n, d)$ requires $\Omega(nd)$ queries. As $D(n/2, d)$ is a subclass of d -bounded degree graphs with n vertices, we conclude that at least $\Omega(\frac{nd}{2})$ queries are needed for this class.

Using theorem 3 we now show that this lower bound can be achieved, modulo a constant factor, by a predetermined algorithm.

Theorem 4 *For a constant d , there exists a predetermined search algorithm for the class of d -bounded degree graphs with n vertices, that requires $24dn$ queries asymptotically.*

Proof: Consider a graph $G = (V, E)$ with $\deg(v) \leq d$ for all $v \in V$, and $|V| = n$. Let μ be the query function, that is $\mu(W) = |E \cap (W \times W)|$ for $W \subseteq V$. Recall that the nodes V are identified with the set of labels $\{1, \dots, n\}$. We associate to G a bipartite graph $G' = (V', V'', E')$, where $V' = V'' = \{1, \dots, n\}$, and $(v_1, v_2) \in E'$ iff $(v_1, v_2) \in E$. Note that $\deg(v) \leq d$ for every $v \in V' \cup V''$. We want to reconstruct graph G by applying Theorem 3 to graph G' . The query function for graph G' is $\mu'(W', W'') = |E' \cap (W' \times W'')|$ for $W' \subseteq V'$, $W'' \subseteq V''$. We now show that a query μ' can be simulated by a constant number of queries μ .

Observe the following properties of μ' :

1. $\mu'(W', W'') = \mu'(W'', W')$
2. $\mu'(W, W) = 2\mu(W)$

3. If $W' \cap W'' = \emptyset$ then $\mu'(W', W'') = \mu(W' \cup W'') - \mu(W') - \mu(W'')$
4. If $W_1 \cap W_2 = \emptyset$ then $\mu'(W_1 \cup W_2, W) = \mu'(W_1, W) + \mu'(W_2, W)$ for any $W \in V''$.

For arbitrary $W_1, W_2 \subseteq V$, let $W = W_1 \cap W_2$. Then

$$\begin{aligned} \mu'(W_1, W_2) &= \mu'((W_1 \setminus W_2) \cup W, (W_2 \setminus W_1) \cup W) = \\ &= \mu'(W_1 \setminus W_2, W_2 \setminus W_1) + \mu'(W_1 \setminus W_2, W) + \mu'(W_2 \setminus W_1, W) + \mu'(W, W). \end{aligned}$$

Using properties 1-4, we obtain

$$\begin{aligned} \mu'(W_1, W_2) &= \\ &= \mu((W_1 \setminus W_2) \cup (W_2 \setminus W_1)) - 2\mu(W_1 \setminus W_2) - 2\mu(W_2 \setminus W_1) + \mu(W_1) + \mu(W_2). \end{aligned}$$

Thus, one query μ' can be simulated by five queries μ . By Theorem 3, graph G' , and therefore G , can be reconstructed through $8nd$ queries $\mu'(W'_i, W''_j)$. The number of corresponding queries μ can be optimized, if queries $\mu(W_i)$, $\mu(W_j)$ are computed once. This gives us $24nd + 4d \log_d n + 2n/\log_d n = 24nd(1 + o(1))$ queries μ . \square

4 Case Studies

Here we consider some particular classes of graphs. For some of them, that are not subclasses of degree bounded graphs, we show that our technique still applies. For others, we show that the multiplicative factor can be improved.

4.1 General and c -Colorable Graphs

The results above assume some knowledge about the class that the unknown graph is drawn from. What can be said when no prior information about the structure of the graph is given, i.e. all graphs are possible? The information-theoretic lower bound for this case is immediate. There are $2^{\frac{n(n-1)}{2}}$ labelled graphs with n vertices and each query can yield up to $1 + n(n-1)/2$ answers. Therefore, any algorithm should make at least $\log_{(1+n(n-1)/2)} 2^{\frac{n(n-1)}{2}} = \Omega(\frac{n^2}{4 \log_2 n})$ queries. Note that since graphs can be represented by $(0, 1)$ -vectors of length $n(n-1)/2$, any non-adaptive algorithm for reconstructing general graphs gives a 2-detecting matrix with $n(n-1)/2$ columns. By Lemma 1, we can then obtain a better lower bound of $2 \cdot \frac{n(n-1)}{2} / \log_2 \frac{n(n-1)}{2} = \Omega(\frac{n^2}{2 \log_2 n})$ for non-adaptive algorithms for reconstructing general graphs. This lower bound can be achieved, within a factor of 4, using the technique of the proof of Theorem 4. Represent $G = (V, E)$ as a bipartite graph $G' = (V_1, V_2, E')$ where each side consists of a copy of vertices of V ($V_1 = V_2 = V$) and $(v_1, v_2) \in E'$ iff $(v_1, v_2) \in E$. For each vertex $v_1 \in V_1$, we can find all its adjacent vertices in V_2 (and then in G) through $\frac{2n}{\log_2 n}$ queries of the form ‘‘How many adjacent vertices does v_1 have in a subset $W \subseteq V_2$?’’. Every such query can be simulated by two queries to the initial graph G . Similar

to the proof of Theorem 4, let μ (resp. μ') denote the query function for graph G (resp. G'). Then $\mu'(\{v_1\}, W) = \mu(W \cup \{v_1\}) - \mu(W \setminus \{v_1\})$. To reconstruct the graph, we find for every vertex i the adjacent vertices among $1, \dots, i-1$. Then the overall complexity is $2 \cdot \sum_{i=2}^n \frac{2i}{\log_2 i} = \frac{2n^2}{\log_2 n} (1 + o(1))$ which is 4 times the lower bound.

Does the knowledge of the graph's chromatic number $c = \chi(G)$ help? Not much, as there are at least $2^{\frac{c-1}{2c}n^2}$ such graphs (divide n vertices into c parts evenly and consider all possible edge combinations between different parts). The information-theoretic lower bound is then $\Omega(\frac{c-1}{4c} \cdot \frac{n^2}{\log n})$, and the algorithm above for the general case is again optimal up to a constant factor.

4.2 h -Edge Colorable Graphs

If the graph is known to be h -edge-colorable, the degree of vertices is bounded by h and by Theorem 4, it can be reconstructed within $O(hn)$ non-adaptive queries. Note that this is asymptotically best possible, as by Vizing's theorem (see [5]), the graphs with the edge chromatic number less than or equal to h contain all the $(h-1)$ -bounded degree graphs.

4.3 Matchings in Bipartite Graphs

Matchings occur in numerous applications and we consider important to present a refinement of the general technique that can be obtained for this class. This refinement is valid for a more general class, namely the 1-bounded one-sided bipartite graphs (see Sect. 3). Let $G = (V, W, E)$ be a bipartite graph, where $|V| = n$, $|W| = m$ and all vertices in V have degree at most 1. According to the proof of Theorem 3, to reconstruct G we need a 1-separating matrix for m objects and 2-detecting matrix for n objects. A 1-separating matrix with m columns has $\log m$ rows (see Sect. 2.1). As for 2-detecting matrix with n columns, $\frac{2n}{\log_2 n} (1 + o(1))$ rows are necessary and sufficient, as it was mentioned in Sect. 2.2. Putting together, G can be reconstructed by $\log_2(m) \frac{2n}{\log_2 n} (1 + o(1))$ non-adaptive queries. Note that the probabilistic proof of Theorem 1 is not involved here, and the queries can be constructed effectively.

In the case of matchings in bipartite graphs we have $n = m$ which gives a non-adaptive algorithm to reconstruct a matching within $2n(1 + o(1))$ queries. This bound is asymptotically optimal within a factor of 2.

4.4 Hamiltonian Cycles and Paths

Let us consider the 2-bounded degree graphs with n vertices. Any such graph is a collection of paths and cycles without common vertices. In particular, this class contains the Hamiltonian cycles and the Hamiltonian paths. The problem of reconstructing Hamiltonian cycles under different models was considered in [10] in connection with its application to genome physical mapping.

Theorem 4 suggests a non-adaptive solution that requires $48n$ queries asymptotically. A better performance can be obtained if we sacrifice the requirement for the algorithm to be fully non-adaptive. Instead, we propose a two-stage algorithm – the first stage does an adaptive “pre-processing” and the second stage reconstructs the graph non-adaptively.

At the first stage, we sort out the vertices into three disjoint independent subsets, that is without adjacent pairs in each subset. As each vertex has at most two neighbours, this sorting can be easily done in at most $2n$ queries by processing the vertices consecutively and testing each vertex against at most two of the already formed subsets.

At the second stage, we reconstruct separately each of the three bipartite 2-degree bounded graphs resulting from the first stage. Again, applying the proof of Theorem 4, we need a 2-separating and a 3-detecting matrices. As noted in Sect. 2.1, a 2-separating matrix with n columns and $2 \log_2 n$ rows can be effectively constructed. On the other hand, by adapting the proof of Theorem 2 to the case $d = 3$, it can be shown that there exists a 3-detecting matrix with n columns and $\frac{4n}{\log n}$ rows. By applying the algorithm of reconstructing bipartite graphs in such a way that the detecting matrix always acts on a smaller part (see proof of Theorem 3), we can reconstruct each bipartite graph in $2 \log_2 n \cdot \frac{4n/2}{\log n/2} = 4n$ queries.

Putting two stages together, this gives an algorithm with $2n + 3 \cdot 4n = 14n$ queries.

5 Remarks

An interesting open question is to give an explicit construction of d -separating matrices with $4 \log_d n$ rows, the existence of which has been proved in Theorem 1 by a probabilistic method. An explicit construction could open a way to the study of computational complexity of reconstructing the unknown graph given a vector of answers. This question is another direction for future research.

References

- [1] Martin Aigner. Search problems on graphs. *Discrete Applied Mathematics*, 14:215–230, 1986.
- [2] Martin Aigner. *Combinatorial Search*. John Wiley & Sons, 1988.
- [3] Noga Alon. Separating matrices. private communication, May 1997.
- [4] Noga Alon and Joel Spencer. *The Probabilistic Method*. Wiley Interscience, New York, 1992.
- [5] Peter J. Cameron. *Combinatorics: Topics, Techniques, Algorithms*. Cambridge University Press, 1994.
- [6] Ding-Zhu Du and Frank K. Hwang. *Combinatorial Group Testing and its applications*, volume 3 of *Series on applied mathematics*. World Scientific, 1993.

- [7] L. Gargano, V. Montuori, G. Setaro, and U. Vaccaro. An improved algorithm for quantitative group testing. *Discrete Applied Mathematics*, 36:299–306, 1992.
- [8] Oded Goldreich, Shafi Goldwasser, and Dana Ron. Property testing and its connection to learning and approximation. In *Proceedings of the 37th Annual Symposium on Foundations of Computer Science*, pages 339–348, 1996.
- [9] Oded Goldreich and Dana Ron. Property testing in bounded degree graphs. In *The 29th Annual ACM Symposium on Theory of Computing*, 1997. to appear. Full version available at <http://theory.lcs.mit.edu/~danar/papers.html>.
- [10] Vladimir Grebinski and Gregory Kucherov. Optimal query bounds for reconstructing a hamiltonian cycle in complete graphs. In *Proceedings of the 5th Israeli Symposium on Theory of Computing and Systems*. IEEE Press, June 1997. to appear.
- [11] Fred H. Hao. The optimal procedures for quantitative group testing. *Discrete Applied Mathematics*, 26:79–86, 1990.
- [12] V. Koubek and J. Rajlich. Combinatorics of separation by binary matrices. *Discrete Mathematics*, 57:203–208, 1985.
- [13] Bernt Lindström. On a combinatorial problem in number theory. *Canad. Math. Bull.*, 8:477–490, 1965.
- [14] Bernt Lindström. On Möbius functions and a problem in combinatorial number theory. *Canad. Math. Bull.*, 14(4):513–516, 1971.
- [15] Bernt Lindström. On B_2 -sequences of vectors. *Journal of Number Theory*, 4:261–265, 1972.
- [16] Bernt Lindström. Determining subsets by unramified experiments. In J.N. Srivastava, editor, *A Survey of Statistical Designs and Linear Models*, pages 407–418. North Holland, Amsterdam, 1975.
- [17] R.L. Rivest and J. Vuillemin. On reconstructing graph properties from adjacency matrices. *Theoretical Computer Science*, 3:371–384, 1976.

Appendix A

In this section we give main ideas of an explicit construction of d -detecting matrix presented in a series of works of B. Lindström, especially in [14].

Consider the columns of a d -detecting matrix as a set of vectors. Their property can be rephrased as follows.

Definition 5.1 *A set of $(0, 1)$ vectors $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n$ is said to be d -detecting iff all the sums $\sum_{i=1}^n \epsilon_i \vec{v}_i$ ($\epsilon_i = 0, 1, \dots, d-1$) are different. \square*

Given n and d , we are interested in a set of d -detecting vectors of *minimal* dimension. Below we outline the asymptotically optimal construction given by Lindström in [14]. The proofs can be found in the original paper.

To introduce the Lindström's construction we recall the definition of Möbius function for partially ordered sets.

Definition 5.2 (Möbius function) *Let (P, \leq) be a finite partially ordered set. The Möbius function $\mu(x, y)$ of P is defined for $x, y \in P$ as follows.*

1. $\mu(x, x) = 1$
2. if $x \not\leq y$ then $\mu(x, y) = 0$
3. if $x < y$ then $\mu(x, y) = -\sum_{x < z < y} \mu(z, y)$

For example, it is known that if P is the Boolean algebra of all subsets of a finite set then

$$\mu(x, y) = (-1)^{|y \setminus x|} \quad \text{if } x \subset y$$

We will use this fact later. Lindström proved the following results:

Theorem 5 ([14]) *Let P be a finite partially ordered set with 0 and a unique last element 1. Let $\mu(x, y)$ be the Möbius function of P . Set $m = \sum_{x \in P} |\mu(x, 1)|$. m is then an even integer. Let n be an arbitrary integer in the interval $0 \leq n \leq m/2$. Then there exists a function $f(x) \in \{0, 1\}$ on P such that*

$$\sum_{0 < x \leq 1} f(x) \mu(x, 1) = -n \cdot \text{sign}(\mu(0, 1)),$$

where $\text{sign}(a) = 1$ if $a \geq 0$ and $\text{sign}(a) = -1$ if $a < 0$.

Theorem 6 ([13]) *Let P be a finite semilattice with Möbius function $\mu(x, y)$. Let $a, b \in P$ and $b \not\leq a$. Let $f(x)$ be defined for all $x \leq a \wedge b$ with values in a commutative ring with unit. Then we have*

$$\sum_{x \leq b} f(x \wedge a) \mu(x, b) = 0 \tag{1}$$

The following theorem is the key-stone of the construction of optimal d -detecting vector set.

Theorem 7 ([14]) *Let (P, \wedge) be a finite semilattice with $m + 1$ elements. Define a partial order on P such that $a \leq b$ iff $a = a \wedge b$. Let θ be the least element in (P, \leq) . Put $m_y = \sum_{x \leq y} |\mu(x, y)|$. Then there exists a d -detecting set containing $\sum_{y > \theta} h_d(m_y/2)$ vectors of dimension m . \square*

Before we give the construction of a d -detecting vector set, we introduce the *detecting capacity* $h_k(x)$.

Definition 5.3 *The detecting capacity $h_k(x)$ is the maximum number h for which there exist integers d_i ($i = 1 \dots h$), $1 \leq d_i \leq x$, such that all the sums $\sum_{i=1}^h \epsilon_i d_i$ ($\epsilon_i = 0 \dots k-1$) are distinct.*

We call a vector $(d_1, \dots, d_h)^T$ on which the maximum is reached a *detecting vector*.

The construction of the d -detecting vector set has the following stages [14]:

1. Given a semilattice $P = (x_1, x_2, \dots, x_m)$ with m elements we consider its multiplication table, an $m \times m$ matrix $(a_{i,j})$, where $a_{i,j} = x_i \wedge x_j$.
2. Consecutively consider each column of the above matrix. Replace each entry of this column by a row-vector of dimension $h_d(m_{x_i}/2)$. Note that the column marked with x_i contains all $x_j \leq x_i$. By theorem 5, for all $k \leq m_{x_i}/2$, there is a function $f_k(x) : P \rightarrow \{0, 1\}$ such that $\sum_{\theta < x \leq x_i} f_k(x) \mu(x, x_i) = -k \cdot \text{sign}(\mu(\theta, x_i))$.
3. For each column x_i , find $h_d(m_{x_i}/2)$ and a corresponding d -detecting vector (d_1, d_2, \dots, d_h) . Then replace the entry at each row x_j by the row-vector $(f_{d_1}(x_i \wedge x_j), f_{d_2}(x_i \wedge x_j), \dots, f_{d_h}(x_i \wedge x_j))$.

To prove that the columns of the obtained matrix form a d -detecting set assume that

$$\sum_{1 \leq i \leq m, 1 \leq j \leq h_d(m_{x_i}/2)} e_{i,j} \vec{v}_{i,j} = \vec{0}, \quad \text{where } e_{i,j} = -k, \dots, 0, \dots, k, \quad (2)$$

We prove (following [14]) that all $e_{i,j} = 0$. If it is not true, then there exists a maximal x_i in (P, \leq) such that $e_{i,j} \neq 0$ for some j . Multiply the v -th row of both side of (2) by $-\mu(x_v, x_i) \cdot \text{sign}(\mu(\theta, x_i))$ and sum up all the rows. It follows from (1) that columns corresponding to $y \not\leq x_i$ sum up to zero. From (2) we further get

$$\sum_{j=1}^h e_{i,j} d_{i,j} = 0,$$

where $h = h_d(m_i/2)$. If some $e_{i,j} \neq 0$ we get a contradiction to the fact that $\{d_{i,j}\}_{j=1}^h$ is detecting.

Now we apply the above construction in order to obtain an upper bound of a d -detecting set of n vectors.

1. Consider the Boolean lattice of subsets of n -element set. So that if $x \subset y$ then $\mu(x, y) = (-1)^{|y \setminus x|}$. It follows that $m_x = 2^{|x|}$.
2. Consider a prefix of the sequence $(1, d, d^2, \dots, d^i, \dots)$. Clearly, this is a d -detecting vector.

3. Now we have $h_d(m) = \lfloor \log_d m \rfloor$, and the detecting capacity of element x is

$$h_d(2^{|x|}/2) = \lfloor \log_d 2^{|x|-1} \rfloor = \lfloor \frac{|x|-1}{\log_2(d)} \rfloor \geq \frac{|x|}{\log_2(d)} - 2.$$

4. This gives a matrix with

$$\sum_{x \subset [n]} \left(\frac{|x|}{\log_2 d} - 2 \right) = \sum_{i=0}^n \binom{n}{i} \left(\frac{i}{\log_2 d} - 2 \right) = \frac{2^n n}{2 \log_2 d} - 2^{n+1}$$

columns.

It follows that there exists a d -detecting set of approximately $\frac{2^n n}{2 \log_2 d}$ vectors of dimension 2^n . We conclude that for large n , there exists a d -detecting set of n vectors of dimension $2 \log(d) \frac{n}{\log n}$ asymptotically.

Appendix B

Here we give a proof of Lemma 1 that an optimal d -detecting matrix has at least $2 \log d \frac{n}{\log n}$ rows asymptotically. The proof is a modification of L.Moser's method described in [16, p. 415].

Let the $\{\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n\}$ be the columns of a d -detecting matrix. We show that these vectors must have dimension at least $2 \log d \frac{n}{\log n}$ asymptotically. Let m be the dimension of \vec{v}_i 's. The idea of the proof is to show that at least half of the vectors of the set $\{\epsilon_1 \vec{v}_1 + \epsilon_2 \vec{v}_2 + \dots + \epsilon_n \vec{v}_n \mid \epsilon_i = 0, \dots, d-1\}$ belongs to an m -dimensional sphere of a "small" radius. This gives an estimation for m .

Consider the uniform probabilistic space $\{(\epsilon_1, \epsilon_2, \dots, \epsilon_n) \mid \epsilon_i = 0, \dots, d-1\} \cong [0, \dots, d-1]^n$. The random variable $\xi = \epsilon_1$ has the expectation $E(\xi) = (d-1)/2$ and the variance $\sigma^2 = \text{Var}(\xi) = \frac{(d+2)d}{12}$. Denote by v_i^j the value of j -th coordinate of \vec{v}_i . Then the random variable $\varsigma_k = \epsilon_1 v_1^k + \dots + \epsilon_n v_n^k$ ($k = 1 \dots m$) is a sum of independent random variables ϵ_i with coefficients 0, 1. It means that

$$\text{Var}(\varsigma_k) \leq \sum_{i=1}^n \text{Var}(\epsilon_i) = n\sigma^2.$$

By definition of variance and the linearity of expectation, we have:

$$E \left(\sum_{k=1}^m (\varsigma_k - \bar{\varsigma}_k)^2 \right) = \sum_{k=1}^m E(\varsigma_k - \bar{\varsigma}_k)^2 \leq m * n\sigma^2.$$

It follows from the Chebyshev inequality,

$$\text{Prob} \left(\sum_{k=1}^m (\varsigma_k - \bar{\varsigma}_k)^2 \leq 2mn\sigma^2 \right) \geq 1/2.$$

Since $(\epsilon_1, \dots, \epsilon_n) \rightarrow \epsilon_1 \vec{v}_1 + \dots + \epsilon_n \vec{v}_n = (\varsigma_1, \dots, \varsigma_m)^T$ is a bijection, at least a half of these sums belong to a sphere with center $(\bar{\varsigma}_1, \dots, \bar{\varsigma}_m)$ and radius $r = (2mn\sigma^2)^{1/2}$. By the volume argument, the number of integer-valued points in a sphere with radius r is less than $(c/m)^{m/2} r^m$ for a constant c . Therefore, $1/2 \cdot d^n \leq (c/m \cdot r^2)^{m/2} = (c/m \cdot 2mn\sigma^2)^{m/2} = (2cn\sigma^2)^{m/2}$. It follows that $m \geq 2 \cdot \frac{n \log d - \log(2)}{\log n + \log(2 \cdot c \cdot (d+2) \cdot d/12)} = 2 \log d \frac{n}{\log(n)} + o(\frac{n}{\log n})$. \square



Unit e de recherche INRIA Lorraine, Technop ole de Nancy-Brabois, Campus scientifique,
615 rue du Jardin Botanique, BP 101, 54600 VILLERS L ES NANCY
Unit e de recherche INRIA Rennes, Irisa, Campus universitaire de Beaulieu, 35042 RENNES Cedex
Unit e de recherche INRIA Rh one-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN
Unit e de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex
Unit e de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

 diteur
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)
ISSN 0249-6399