

A Characterization of Easily Testable Induced Subgraphs

Extended Abstract

Noga Alon ^{*} Asaf Shapira [†]

September 10, 2003

Abstract

Let H be a fixed graph on h vertices. We say that a graph G is *induced H -free* if it does not contain any *induced* copy of H . Let G be a graph on n vertices and suppose that at least ϵn^2 edges have to be added to or removed from it in order to make it induced H -free. It was shown in [5] that in this case G contains at least $f(\epsilon, h)n^h$ induced copies of H , where $1/f(\epsilon, h)$ is an extremely fast growing function in $1/\epsilon$, that is independent of n . As a consequence it follows that for every H , testing induced H -freeness with one-sided error has query complexity independent of n . A natural question, raised by the first author in [1], is to decide for which graphs H the function $1/f(\epsilon, H)$ can be bounded from above by a polynomial in $1/\epsilon$. An equivalent question is, for which graphs H , can one design a one-sided error property tester for testing induced H -freeness, whose query complexity is polynomial in $1/\epsilon$. We answer this question by showing that, quite surprisingly, for any graph other than the paths of lengths 1,2 and 3, the cycle of length 4, and their complements, no such property tester exists. We further show that a similar result also applies to the case of directed graphs, thus answering a question raised by the authors in [9]. We finally show that the same results hold even in the case of two-sided error property testers. The proofs combine combinatorial, graph theoretic and probabilistic arguments with results from additive number theory.

1 Preliminaries

1.1 Definitions

All graphs and directed graphs (=digraphs) considered here are finite and have no loops and no parallel edges. Let P be a property of graphs, that is, a family of graphs closed under isomorphism.

^{*}Schools of Mathematics and Computer Science, Raymond and Beverly Sackler Faculty of Exact Sciences, Tel Aviv University, Tel Aviv, Israel. Email: noga@math.tau.ac.il. Research supported in part by a USA-Israeli BSF grant, by the Israel Science Foundation and by the Hermann Minkowski Minerva Center for Geometry at Tel Aviv University.

[†]School of Computer Science, Raymond and Beverly Sackler Faculty of Exact Sciences, Tel Aviv University, Tel Aviv, Israel. Email: asafico@math.tau.ac.il. This work forms part of the author's Ph.D. Thesis.

A graph G with n vertices is ϵ -far from satisfying P if no graph \tilde{G} with the same vertex set, which differs from G in at most ϵn^2 places, (i.e., can be constructed from G by adding and removing at most ϵn^2 edges), satisfies P . An ϵ -tester, or *property tester*, for P is a randomized algorithm which, given the quantity n and the ability to make queries whether a desired pair of vertices of an input graph G with n vertices are adjacent or not, distinguishes with probability at least $\frac{2}{3}$ between the case of G satisfying P and the case of G being ϵ -far from satisfying P . Such an ϵ -tester is a *one-sided* ϵ -tester if when G satisfies P the ϵ -tester determines that this is the case (with probability 1). The ϵ -tester is a *two-sided* ϵ -tester if it may determine that G does not satisfy P even if G satisfies it. Obviously, the probability $\frac{2}{3}$ appearing above can be replaced by any constant smaller than 1, by repeating the algorithm an appropriate number of times.

The property P is called *strongly-testable*, if for every fixed $\epsilon > 0$ there exists a one-sided ϵ -tester for P whose total number of queries is bounded only by a function of ϵ , which is independent of the size of the input graph. This means that the running time of the algorithm is also bounded by a function of ϵ only, and is independent of the input size.

In what follows we denote by P_2, P_3 and P_4 the paths of lengths 1, 2 and 3 (which have 2, 3 and 4 vertices respectively), and by C_4 , the cycle of length 4. We measure query-complexity by the number of vertices sampled, assuming we always examine all edges spanned by them. For a fixed graph H , let P_H^* denote the property of being induced H -free. Therefore, G satisfies P_H^* if and only if it contains no induced subgraph isomorphic to H . We define P_H to be the property of being (not necessarily induced) H -free. Therefore, G satisfies P_H if and only if it contains no copy of H . Thus, for example, for $H = C_4$, any clique of size at least 4 satisfies P_H^* , but does not satisfy P_H .

1.2 Related work

The general notion of property testing was first formulated by Rubinfeld and Sudan [29], who were motivated mainly by its connection to the study of program checking. The study of the notion of testability for combinatorial objects, and mainly for labelled graphs, was introduced by Goldreich, Goldwasser and Ron [23], who showed that all graph properties describable by the existence of a partition of a certain type, and among them k -colorability, have efficient ϵ -testers. The fact that k -colorability is strongly testable is, in fact, implicitly proven already in [12] for $k = 2$ and in [27] (see also [2]) for general k , using the Regularity Lemma of Szemerédi [30], but in the context of property testing it is first studied in [23], where far more efficient algorithms are described. These have been further improved in [7]. The notion of property testing has been investigated in other contexts as well, including the context of regular languages, [6], functions [21], [8], [3], [19], [20], computational geometry [16], [4], graph and hypergraph coloring [15], [8], [11] and other contexts. See [22], [28] and [18] for surveys on the topic.

2 The Main Results

2.1 Background

In [5] it is shown that every first order graph property without a quantifier alternation of type “ $\forall\exists$ ” has ϵ -testers whose query complexity is independent of the size of the input graph. It follows from the main result of [5] that for every fixed H , the property P_H^* is strongly testable. Although the query complexity is independent of n , it has a huge dependency on $1/\epsilon$ (the third function in the Ackerman Hierarchy, which is a tower of towers of exponents). In [2], it was shown, using Szemerédi’s Regularity Lemma, that for every fixed H , the property P_H is also strongly testable. This result was generalized to the case of directed graphs (=digraphs) in [9], by first proving a directed version of the regularity lemma. In the above two cases the query complexity is also huge, a tower of 2’s of height polynomial in $1/\epsilon$. For some graphs, however, there are obviously much more efficient property testers than the ones guaranteed by the above general results. For example for the case of H being an edge, there is obviously a one-sided error property tester for both P_H and P_H^* , whose query complexity is $\Theta(1/\epsilon)$. A natural question is therefore, to decide for which graphs H can one design a one-sided error property tester for P_H or P_H^* , whose query complexity would be bounded by a polynomial of $1/\epsilon$. We call a property P *easily testable* if there is a one-sided error property tester for P whose query complexity is polynomial in $1/\epsilon$.

In [1] it is shown that for an undirected graph H , P_H is easily testable if and only if H is bipartite. The authors of [9] obtain a precise characterization of all digraphs H for which P_H is easily testable. As is evident from this characterization, recognizing these digraphs is rather difficult. Indeed, it is shown in [9] that deciding whether for a directed graphs H , P_H is easily testable, is *NP*-complete. The next natural steps, suggested in [1] and [9], are therefore to give characterizations of the graphs and digraphs H for which P_H^* is easily testable. In this paper we give such characterizations.

2.2 The new results

Our first new result here is the following:

Theorem 1 *Let H be a fixed undirected graph other than P_2, P_3, P_4, C_4 and their complements. Then, there exists a constant $c = c(H) > 0$ such that the query-complexity of any one-sided error ϵ -tester for P_H^* is at least*

$$\left(\frac{1}{\epsilon}\right)^{c \log(1/\epsilon)}.$$

As P_2 -freeness can obviously be tested with query complexity $\Theta(1/\epsilon)$, the following theorem, together with the above theorem, supplies a complete characterization for the graphs H for which P_H^* is easily testable, except for the case of P_4, C_4 and its complement (the complement of P_4 is also P_4).

Theorem 2 *There is a one-sided error property tester for testing P_3 -freeness, with query complexity*

$$O(\log(1/\epsilon)/\epsilon).$$

We also prove the following theorem, which is analogous to Theorem 1, only with respect to directed graphs.

Theorem 3 *Let H be a fixed directed graph on at least 5 vertices. Then, there exists a constant $c = c(H) > 0$ such that the query-complexity of any one-sided error ϵ -tester for P_H^* is at least*

$$\left(\frac{1}{\epsilon}\right)^{c \log(1/\epsilon)}.$$

We can actually prove a super-polynomial (in $1/\epsilon$) lower bound for the query complexity of P_H^* for some of the directed graphs H on at most 4 vertices as well.

We finally show that Theorems 1 and 3 can also be extended to the cases of two-sided error property testers.

Theorem 4 *All the lower bounds of Theorems 1 and 3 hold for two-sided error property testers as well.*

2.3 Organization

The (short) proof of Theorem 2 appears in section 3. The lower bound proved by Theorem 1 is established in section 4. To prove this result we have to construct, for any graph H (other than the ones mentioned in the theorem) and any small $\epsilon > 0$, a graph G which is ϵ -far from being induced H -free and yet contains relatively few induced copies of H . The proof of this part, described in Section 4, uses the approaches of [1] and [9], but requires several additional ideas. It applies certain constructions in additive number theory, based on (simple variants of) the construction of Behrend [10] of dense subsets of the first n integers without three-term arithmetic progressions. The proof of Theorem 3 also appears in section 4. In Section 5 we sketch the proof of Theorem 4 which extends the lower bounds of Theorems 1 and 3 to the more general cases of two-sided error property-testers. Due to space limitations, the complete proof appears in appendix 7.2. The final section, Section 6, contains some concluding remarks and open problems.

Throughout this extended abstract we assume, whenever this is needed, that the number of vertices n of the graphs or digraph G is sufficiently large, and that the error parameter ϵ is sufficiently small. In order to simplify the presentation, we omit all floor and ceiling signs whenever these are not crucial, and make no attempt to optimize the absolute constants. In order to make the presentation simple, from this part on, when we refer to a *copy* of H in G , we always refer to an *induced copy* of H in G . Also, when we speak of the property of being H -free, we mean the property of being *induced H -free*.

3 Easily Testable Graphs

In this section we describe the proof of Theorem 2. The algorithm simply picks a random subset of $O(\log(1/\epsilon)/\epsilon)$ vertices, and checks if there is an induced copy of P_3 spanned by the set. If G is P_3 -free, the algorithm clearly always answers correctly. We therefore only have to show that if G is ϵ -far from being P_3 -free, the algorithm finds an induced copy of P_3 with probability at least $2/3$. In the proof we use the well known observation that a graph G is P_3 -free, if and only if it is a disjoint union of cliques.

Proof of Theorem 2 Let T denote the set of vertices of G whose degree is at least $\frac{\epsilon}{4}n$, and let W be any set of vertices that contains all the vertices of T but at most $\frac{\epsilon}{4}n$ of them. We claim that if G is ϵ -far from being P_3 -free, then the induced subgraph of G on W is at least $\frac{\epsilon}{2}$ -far from being P_3 -free. Assume, this is not the case. Then we can make less than $\frac{\epsilon}{2}n^2$ changes within W and get a graph that contains no induced copy of P_3 within W . Finally, we remove all the edges that touch a vertex not in $T \cup W$ (there are at most $\frac{\epsilon}{4}n^2$ such edges), and any edge that touches a vertex in $T \setminus W$ (there are at most $\frac{\epsilon}{4}n^2$ such edges), and thus get a graph that is P_3 -free. As altogether we make less than ϵn^2 changes in G , this contradicts the assumption that G is ϵ -far from being T -free.

Let A be a randomly chosen subset of size $8 \log(1/\epsilon)/\epsilon$, and consider a vertex $v \in T$. The probability that A does not contain any neighbor of v is at most $(1 - \frac{\epsilon}{4})^{8 \log(1/\epsilon)/\epsilon} \leq \epsilon^2 \leq \epsilon/32$. As T is of size at most n , it follows that the expected number of vertices that belong to T and have no neighbor in A , is at most $\frac{\epsilon}{32}n$. By Markov's inequality, with probability at least $7/8$, the number of these vertices is at most $\frac{\epsilon}{4}n$.

Assume we were successful in choosing a set A such that all but at most $\frac{\epsilon}{4}n$ members of T have a neighbor in A . If A contains an induced copy of P_3 we are done, otherwise there is a unique partition of A into cliques C_1, \dots, C_r . If a vertex $v \notin A$ has a neighbor $u \in A$ that belongs to C_i , it follows that if G can be partitioned into cliques, D_1, \dots, D_k , where for $1 \leq i \leq r$, $C_i \subseteq D_i$, then v must belong to D_i . For each vertex $v \notin A$ that has a neighbor $u \in C_i \subseteq A$, assign v the number i . If u has neighbors in A that belong to different cliques, then pick any of these numbers. This numbering induces a partition of all the vertices of G that have a neighbor in A into r cliques. As G is by assumption ϵ -far from being P_3 -free it follows from the discussion at the beginning of the proof, that there are at least $\frac{\epsilon}{2}n^2$ pairs of vertices $u, v \in V(G) \setminus A$, such that either u and v should belong to the same clique, but u and v are not connected, or u and v should belong to different cliques, but u and v are connected. It follows that choosing a set B of $6/\epsilon$ randomly chosen pairs of vertices finds such a violating pair with probability at least $\frac{7}{8}$. The probability of A failing to satisfy the required property is at most $\frac{1}{8}$, and the same applies also for B . It follows that with probability at least $\frac{3}{4}$ the induced subgraph on $A \cup B$ is not P_3 -free. As $|A| + |B| = O(\log(1/\epsilon)/\epsilon)$ the proof is complete. ■

4 Hard to Test Graphs and Digraphs

In this section we give the proofs of Theorems 1 and 3. The approach uses a construction in additive number theory, which uses the technique of Behrend [10], used to construct dense sets of integers with no three-term arithmetic progressions. A set $X \subseteq M = \{1, 2, \dots, m\}$ is called *h-sum-free* if for every three positive integers $a, b, c \leq h$ such that $a + b = c$, if $x, y, z \in X$ satisfy the equation $ax + by = cz$ then $x = y = z$. That is, whenever $a + b = c$, and $a, b, c \leq h$, the only solution to the equation that uses values from X , is one of the $|X|$ trivial solutions. We need the following lemma, whose proof appears in Appendix 7.1 due to space limitations:

Lemma 4.1 *For every positive integer m , there exists an h -sum-free subset $X \subset M = \{1, 2, \dots, m\}$ of size at least*

$$|X| \geq \frac{m}{e^{10\sqrt{\log h \log m}}} \quad (1)$$

We proceed with the proofs of Theorems 1 and 3. It is convenient to start the discussion with directed graphs and then obtain the results for undirected graphs as a special case, (as they can be viewed as symmetric digraphs).

An *s-blow-up* of a digraph $H = (V(H), E(H))$ on h vertices is the digraph obtained from H by replacing each vertex $v_i \in V(H)$ by an independent set I_i of size s , and each directed edge $(v_i, v_j) \in E(H)$, by a complete bipartite directed subgraph whose vertex classes are I_i and I_j , and whose edges are directed from I_i to I_j . Note that if we take an s -blow-up of H , we get a graph on sh vertices that contains at least s^h induced copies of H , where each vertex of the copy belongs to a different blow-up of a vertex from H (simply pick one vertex from each independent set). We call these copies the *special copies* of the blow-up. As each pair of vertices in the blow-up is contained in at most s^{h-2} special copies of H , it follows that adding or removing an edge from the graph can destroy at most s^{h-2} copies of H . It follows that one must add or remove at least $s^h/s^{h-2} = s^2$ edges from the blow-up in order to destroy all its special copies of H .

For the proofs of Theorems 1 and 3, we will need the following lemma, in which a triangle in a digraph is simply three vertices u, v, w , such that there is at least one edge between each of the three pairs.

Lemma 4.2 *For every fixed directed graph $H = (V(H), E(H))$ on h vertices, that contains at least one triangle, there is a constant $c = c(H) > 0$, such that for every positive $\epsilon < \epsilon_0(H)$ and every integer $n > n_0(\epsilon)$, there is a digraph G on n vertices which is ϵ -far from being induced H -free, and yet contains at most $\epsilon^{c \log(1/\epsilon)} n^h$ induced copies of H .*

Proof: Given a small $\epsilon > 0$, let m be the largest integer satisfying

$$\epsilon \leq \frac{1}{h^4 e^{10\sqrt{\log m \log h}}}. \quad (2)$$

It is easy to check that this m satisfies

$$m \geq \left(\frac{1}{\epsilon}\right)^{c \log(1/\epsilon)} \quad (3)$$

for an appropriate $c = c(h) > 0$. Let $X \subset \{1, 2, \dots, m\}$ be as in Lemma 4.1. Call the vertices of H v_1, \dots, v_h , and let V_1, V_2, \dots, V_h be pairwise disjoint sets of vertices, where $|V_i| = im$ and we denote the vertices of V_i by $\{1, 2, \dots, im\}$, where, with a slight abuse of notation, we think on the sets V_i as being pairwise disjoint. For each j , $1 \leq j \leq m$, for each $x \in X$ and for each directed edge (v_p, v_q) of H , let $j + (p-1)x \in V_p$ have an outgoing edge pointed to $j + (q-1)x \in V_q$. In other words, for each $1 \leq j \leq m$ and $x \in X$, the graph F contains a copy of H , denoted by $H_{j,x}$, which is spanned by the vertices $j, j+x, j+2x, \dots, j+(h-1)x$. Note that each of these $m|X|$ copies of H is spanned by a set of vertices that forms an arithmetic progression whose first element is j and whose difference is x . A crucial implication is that F contains $m|X|$ copies of H , such that each pair of copies have at most one common vertex. In what follows we call these $m|X|$ copies of H in F , the *essential* copies of H in F . Finally, define

$$s = \left\lfloor \frac{n}{|V(H)|} \right\rfloor = \left\lfloor \frac{2n}{h(h+1)m} \right\rfloor$$

and let G be the s -blow-up of F (together with some isolated vertices, if needed, to make sure that the number of vertices is precisely n).

We now turn to show that G is indeed ϵ -far from being H -free. Consider two essential copies of H in F , H_1 and H_2 . As was noted above, H_1 and H_2 share at most one vertex v_i in F . It follows that their corresponding blow-ups in G will share at most one common independent set I_i , which replaces the vertex v_i from F . Now, consider the blow-ups of H_1 and H_2 in G , denoted T_1 and T_2 . As T_1 and T_2 share at most one common independent set, we conclude that adding or removing an edge from G , can either destroy special copies of H that belong to T_1 , or special copies of H that belong to T_2 (or not destroy any copies at all). As was explained above, in order to destroy all the special copies of an s -blow-up of H , one needs to add or remove at least s^2 edges from the blow-up. As G contains $m|X|$ blow-ups of essential copies of H , we conclude that one has to add or delete at least

$$s^2 m |X| = \frac{4m|X|n^2}{h^2(h+1)^2 m^2} \geq \frac{|X|n^2}{h^4 m} \geq \frac{n^2}{h^4 e^{10\sqrt{\log m \log h}}} \geq \epsilon n^2 \quad (4)$$

edges in order to make G H -free. The second inequality follows from the lower bound on $|X|$ guaranteed by Lemma 4.1, and the third from (2). It follows that G is indeed ϵ -far from being H -free.

To complete the proof, we have to show that G contains at most $\epsilon^{c \log(1/\epsilon)} n^h$ induced copies of H . Note, that as H contains at least one triangle, and each triangle belongs to at most $\binom{n}{h-3} \leq n^{h-3}$ copies of H , it is enough to show that G contains at most $\epsilon^{c \log(1/\epsilon)} n^3$ triangles. Consider a partition of the vertices of G into h subsets U_1, \dots, U_h , where U_i contains the im independent sets that resulted

from the blow-ups of the im vertices that belonged to V_i in F . Notice that if we show that the induced subgraph of G on any three of the subsets U_1, \dots, U_h contains at most $\epsilon^{c' \log(1/\epsilon)} n^3$ triangles, then the total number of triangles in G is at most $\binom{h}{3} \epsilon^{c' \log(1/\epsilon)} n^3$, which is still at most $\epsilon^{c \log(1/\epsilon)} n^3$.

Fix any three subsets U_i, U_j, U_k such that $i < j < k$, and recall that they are the result of an s -blow-up of V_i, V_j, V_k . As there are no edges within these sets any triangle spanned by them must have exactly one vertex in each set. Note, that if the sets span a triangle whose vertices belong to the independent sets $I_x \subseteq U_i, I_y \subseteq U_j, I_z \subseteq U_k$, then the vertices $x \in V_i, y \in V_j, z \in V_k$ in F must also span a triangle. It follows that the number of triangles spanned by U_i, U_j, U_k is exactly s^3 times the number of triangles spanned by V_i, V_j, V_k .

If the vertices v_i, v_j, v_k , do not span a triangle in H , then by the definition of F , V_i, V_j, V_k do not span a triangle, and so do U_i, U_j, U_k in G , and we are done. If v_i, v_j, v_k span a triangle in H , then by definition of F for any triangle spanned by V_i, V_j, V_k , there are $x, y \in X$ and $1 \leq t \leq im$, such that the three vertices of this triangle are $t \in V_i, t + (j - i)x \in V_j, t + (j - i)x + (k - j)y \in V_k$. As this is a triangle, there must also be some $z \in X$ such that $t + (j - i)x + (k - j)y = t + (k - i)z$. We conclude that the following equation in positive coefficients, whose values are at most h , holds

$$(j - i)x + (k - j)y = (k - i)z.$$

As X is h -sum free, it follows that $x = y = z$, hence V_i, V_j, V_k span precisely $m|X|$ triangles of the form $t + (i - 1)x \in V_i, t + (j - 1)x \in V_j, t + (k - 1)x \in V_k$, for every $1 \leq t \leq m$ and $x \in X$. We conclude that U_i, U_j, U_k span at most $m|X|s^3 < m^2(n/m)^3 \leq n^3/m$ triangles. As by (3), $m \geq (1/\epsilon)^{c \log(1/\epsilon)}$, the proof is complete. \blacksquare

The proofs of Theorems 1 and 3 now follow easily from the above Lemma.

Proof of Theorem 1: Let H be a fixed graph on h vertices. A simple yet crucial observation is that for every graph H testing P_H^* is equivalent to testing $P_{\overline{H}}^*$, where \overline{H} is the complement of H . Note, that this relation does not hold for testing P_H . It follows that in order to prove a lower bound for testing P_H^* , we may prove a lower bound for testing $P_{\overline{H}}^*$.

Given a one-sided error ϵ -tester for testing H -freeness we may assume, without loss of generality, that it queries about all pairs of a randomly chosen set of vertices (otherwise, as explained in [5], every time the algorithm queries about a vertex pair we make it query also about all pairs containing a vertex of the new pair and a vertex from previous queries. This may only square the number of queries. See also [24] for a more detailed proof of this statement.) As the algorithm is a one-sided-error algorithm, it can report that G is not H -free only if it finds a copy of H in it. Observe that if the tester picks a random subset of x vertices, and an input graph contains only $\epsilon^{c \log(1/\epsilon)} n^h$ copies of H , then the expected number of copies of H spanned by x is at most $x^h \epsilon^{c \log(1/\epsilon)}$, which is far smaller than 1 unless x exceeds $(1/\epsilon)^{c' \log(1/\epsilon)}$ for some $c' = c'(H) > 0$. It follows by Markov's inequality that the tester finds a copy of H with negligible probability. It is therefore enough to show that for

any undirected graph H , other than P_2, P_3, P_4, C_4 and their complements, there is a graph G on n vertices which is ϵ -far from being H -free, yet contains only $\epsilon^{c \log(1/\epsilon)} n^h$ copies of H .

If $h \geq 6$, then it follows from the simplest result in Ramsey Theory (c.f., e.g., [25], page 1) that either H or \overline{H} must contain a triangle. Hence, assuming that H contains a triangle, we can use Lemma 4.2 to construct a graph G on n vertices which is ϵ -far from being H -free and yet contains at most $\epsilon^{c \log(1/\epsilon)} n^h$ copies of H . For $h = 5$, the only graph H , such that neither H nor \overline{H} contain a triangle is C_5 (the cycle of length 5, whose complement is also C_5). In this case we can use the fact that C_5 is the core of itself (see [1], [9] for the definition of a core) to prove that $P_{C_5}^*$ is not easily testable. Due to space limitations, we omit the detailed argument. As for $h = 2, 3, 4$ the only graphs H for which H and \overline{H} are triangle-free are P_2, P_3, P_4, C_4 and their complements, the proof is complete. ■

Proof of Theorem 3: The proof is similar to the proof of Theorem 1. One only has to note again that for every directed graph H , on at least 6 vertices, either H or \overline{H} contain a triangle, and that the only digraph on 5 vertices which does not have this property is the directed graph obtained from C_5 , by replacing each undirected edge with two directed edges. This case can be handled directly, as it is essentially equivalent to the undirected case of $H = C_5$. In fact, using some of the methods of [9] we can show that for some of the digraphs H of size 3 and 4, P_H^* is not easily testable. We postpone the details to the full version of the paper. ■

5 Two-Sided Error Property-Testers

For the proof of Theorem 4 we apply Yao's principle [31], by constructing, for every fixed graph H , for which a lower bound was established in Theorems 1 and 3, two distributions D_1 and D_2 , where D_1 consists of graphs which are ϵ -far from being H -free with probability $1 - o(1)$ (where the $o(1)$ term tends to 0 as ϵ tends to zero), while D_2 consists of graphs which are H -free. We then show that any deterministic algorithm, which makes a small number of queries (adaptively) cannot distinguish with non-negligible probability between D_1 and D_2 . We prove Theorem 4 for the case of directed graphs, as it is clear that the case of undirected graphs follows as a special case. For the case of H being the graph obtained from C_5 by replacing each edge by a cycle of length 2, we can use the fact that this graph is the core of itself (see [1], [9] for the definition of a core) to prove that $P_{C_5}^*$ is not easily testable, even with two-sided error. Due to space limitations, we omit the detailed argument. We thus assume that H is a graph on at least 6 vertices. As in the proofs of Theorems 1 and 3, testing P_H^* with two-sided error has the same query complexity as testing $P_{\overline{H}}^*$, thus we assume that H contains at least one triangle. Due to space limitations, the complete proof appears in appendix 7.2.

6 Concluding Remarks and Open Problems

- As in the case of P_H , there is a huge gap between the general upper bounds for testing P_H^* that were established in [5], and the lower bounds in this paper. It would be very interesting, and probably challenging, to improve any of these bounds. Even in the seemingly simplest case of H being a triangle, we do not know how to improve these bounds.
- Another interesting open problem is to complete the characterizations of easily testable properties P_H^* by solving the cases of $H = P_4, C_4$, and by classifying all the directed graphs on 3 and 4 vertices. The case of testing P_4 -freeness seems the simplest one to resolve, since there are known structural results, that characterize P_4 -free graphs. These graphs are also known as Complement Reducible graphs, or Cographs for short, and they are precisely the graphs formed from a single vertex under the closure of the operations of union and complement, see [14], [26]. Cographs have a unique tree representation called a cotree. It might be possible to use this characterization, and the unique tree representation in order to design an efficient tester for P_4 -freeness. For some of the directed graphs H on 3 or 4 vertices it is also possible to decide if P_H^* is easily testable by combining our techniques here with the methods in [9], but a few cases remain. A detailed discussion is postponed to the full version of the paper.
- There is an interesting possible connection between the problem of graph isomorphism and testing P_H^* . It is known (see [13]) that for the graphs P_2, P_3, P_4 and C_4 , the graph isomorphism problem can be solved in polynomial time for H -free graphs. Moreover, for any other H , any instance of the graph isomorphism problem can be reduced to an instance that is H -free. Thus, in some sense, the problem on H -free graphs, for H other than P_1, P_2, P_3 and C_4 , is *isomorphism hard*. It might be interesting to understand if this connection is indeed meaningful.

References

- [1] N. Alon, Testing subgraphs in large graphs, Proc. 42nd IEEE FOCS, IEEE (2001), 434-441.
- [2] N. Alon, R. A. Duke, H. Lefmann, V. Rödl and R. Yuster, The algorithmic aspects of the Regularity Lemma, Proc. 33rd IEEE FOCS, Pittsburgh, IEEE (1992), 473-481. Also: J. of Algorithms 16 (1994), 80-109.
- [3] N. Alon, W. F. de la Vega, R. Kannan and M. Karpinski, Random Sampling and Approximation of MAX-CSP Problems, Proc. of the 34th ACM STOC, ACM Press (2002), 232-239.
- [4] N. Alon, S. Dar, M. Parnas and D. Ron, Testing of clustering, Proc. 41st IEEE FOCS, IEEE (2000), 240-250.

- [5] N. Alon, E. Fischer, M. Krivelevich and M. Szegedy, Efficient testing of large graphs, Proc. of 40th FOCS, New York, NY, IEEE (1999), 656–666. Also: *Combinatorica* 20 (2000), 451-476.
- [6] N. Alon, M. Krivelevich, I. Newman and M. Szegedy, Regular languages are testable with a constant number of queries, Proc. 40th FOCS, New York, NY, IEEE (1999), 645–655. Also: *SIAM J. on Computing* 30 (2001), 1842-1862.
- [7] N. Alon and M. Krivelevich, Testing k-colorability, *SIAM J. Discrete Math.*, 15 (2002), 211-227.
- [8] N. Alon and A. Shapira, Testing satisfiability, Proc. of the 13th Annual ACM-SIAM SODA, ACM Press (2002), 645-654. Also: *Journal of Algorithms*, to appear.
- [9] N. Alon and A. Shapira, Testing subgraphs in directed graphs, Proc. of the 35th Annual Symp. on Theory of Computing (STOC), San Diego, California, 2003, to appear.
- [10] F. A. Behrend, On sets of integers which contain no three terms in arithmetic progression, *Proc. National Academy of Sciences USA* 32 (1946), 331–332.
- [11] A. Bogdanov, K. Obata and L. Trevisan, A Lower Bound for Testing 3-Colorability in Bounded-degree Graphs, Proc. 43rd IEEE FOCS, IEEE (2002), 93-102.
- [12] B. Bollobás, P. Erdős, M. Simonovits and E. Szemerédi, Extremal graphs without large forbidden subgraphs, *Annals of Discrete Mathematics* 3 (1978), 29–41.
- [13] D. G. Corneil, H. Lerchs, L. Stewart Burlingham, Complement Reducible Graphs, *Discrete Applied Mathematics* 3 (1981), 163–174.
- [14] D. G. Corneil, Y. Perl and L. K. Stewart, A linear recognition algorithm for cographs, *SIAM J. Comput.* 14 (1985), 926–934.
- [15] A. Czumaj and C. Sohler, Testing hypergraph coloring, Proc. of ICALP 2001, 493-505.
- [16] A. Czumaj and C. Sohler, Property testing in computational geometry, Proceedings of the 8th Annual European Symposium on Algorithms (2000), 155–166.
- [17] P. Erdős, P. On extremal problems of graphs and generalized graphs. *Israel J. Math.* 2 1964 183-190.
- [18] E. Fischer, The art of uninformed decisions: A primer to property testing, *The Computational Complexity Column of The Bulletin of the European Association for Theoretical Computer Science* 75 (2001), 97-126.
- [19] E. Fischer, G. Kindler, D. Ron, S. Safra, and A. Samorodnitsky, Testing juntas, Proc. of The 43rd FOCS (2002), 103-112.

- [20] E. Fischer, E. Lehman, I. Newman, S. Raskhodnikova, R. Rubinfeld and A. Samorodnitsky, Monotonicity testing over general poset domains, Proc. of The 34th STOC (2002), 474-483.
- [21] A. Frieze and R. Kannan, Quick approximation to matrices and applications, *Combinatorica* 19 (1999), 175-220.
- [22] O. Goldreich, Combinatorial property testing - a survey, In: Randomization Methods in Algorithm Design (P. Pardalos, S. Rajasekaran and J. Rolim eds.), AMS-DIMACS (1998), 45-60.
- [23] O. Goldreich, S. Goldwasser and D. Ron, Property testing and its connection to learning and approximation, *Proceedings of the 37th Annual IEEE FOCS* (1996), 339-348. Also: *Journal of the ACM* 45 (1998), 653-750.
- [24] O. Goldreich and L. Trevisan, Three theorems regarding testing graph properties, Proc. 42nd IEEE FOCS, IEEE (2001), 460-469.
- [25] R. L. Graham, B. L. Rothschild and J. H. Spencer, *Ramsey Theory*, Second Edition, Wiley, New York, 1990.
- [26] T. A. McKee and F.R. McMorris, **Topics in Intersection Graph Theory**, SIAM, Philadelphia, PA, 1999.
- [27] V. Rödl and R. Duke, On graphs with small subgraphs of large chromatic number, *Graphs and Combinatorics* 1 (1985), 91-96.
- [28] D. Ron, Property testing, in: P. M. Pardalos, S. Rajasekaran, J. Reif and J. D. P. Rolim, editors, *Handbook of Randomized Computing*, Vol. II, Kluwer Academic Publishers, 2001, 597-649.
- [29] R. Rubinfeld and M. Sudan, Robust characterization of polynomials with applications to program testing, *SIAM J. on Computing* 25 (1996), 252-271.
- [30] E. Szemerédi, Regular partitions of graphs, In: *Proc. Colloque Inter. CNRS* (J. C. Bermond, J. C. Fournier, M. Las Vergnas and D. Sotteau, eds.), 1978, 399-401.
- [31] A. C. Yao, Probabilistic computation, towards a unified measure of complexity. Proc. of the 18th IEEE FOCS (1977), 222-227.

7 Appendix

7.1 Proof of Lemma 4.1

The main idea of the proof is to construct the set X by representing a set of numbers in base d , and then choosing d in order to maximize the size of X . Let d be an integer (to be chosen later) and define

$$X = \left\{ \sum_{i=0}^k x_i d^i \mid x_i < \frac{d}{h} \ (0 \leq i \leq k) \ \wedge \ \sum_{i=0}^k x_i^2 = B \right\},$$

where $k = \lfloor \log m / \log d \rfloor - 1$ and B is chosen to maximize the cardinality of X .

Assume $x, y, z \in X$ satisfy the equation $ax + by = cz$, where $a, b, c \leq h$ are positive integers, $a + b = c$, and

$$x = \sum_{i=0}^k x_i d^i, \quad y = \sum_{i=0}^k y_i d^i, \quad z = \sum_{i=0}^k z_i d^i.$$

As $x_i, y_i, z_i < d/h$, and $a + b = c$ we have for every i , $0 \leq i \leq k$

$$ax_i + by_i = cz_i.$$

By the convexity of the function $f(z) = z^2$ this implies that

$$ax_i^2 + by_i^2 \geq cz_i^2,$$

and the inequality is strict unless all 3 numbers x_i, y_i and z_i are equal. However, if for some i the inequality is strict, we have

$$a \sum_{i=0}^k x_i^2 + b \sum_{i=0}^k y_i^2 > c \sum_{i=0}^k z_i^2$$

which is impossible as by definition of X

$$\sum_{i=0}^k x_i^2 = \sum_{i=0}^k y_i^2 = \sum_{i=0}^k z_i^2 = B.$$

Thus, $x_i = y_i = z_i$ for all i , and X has no nontrivial solution to the above equation.

We now aim at giving an upper bound for $m/|X|$. We lose a factor of d^2 of the numbers $\{1, \dots, m\}$, due to taking $k = \lfloor \log m / \log d \rfloor - 1$. As we restrict $x_i < d/h$, we lose a factor of h per digit, for a total of h^{k+1} . As $x_i < d/h$, we have $0 \leq \sum_{i=0}^k x_i^2 \leq (k+1)d^2/h^2$. As we chose B to maximize $|X|$, we can not lose more than the average, hence, this is a factor of at most $(k+1)d^2/h^2$. We conclude that

$$|X| \geq \frac{m}{d^2 h^{k+1} (k+1) \frac{d^2}{h^2}}.$$

Taking $d = \lfloor e^{\sqrt{\log m \log h}} \rfloor$, we conclude (with room to spare) that

$$|X| \geq \frac{m}{e^{10\sqrt{\log m \log h}}}. \quad \blacksquare$$

7.2 Proof of Theorem 4

Let H be a fixed digraph which contains at least one triangle. Given n and ϵ , let X , m and the sets V_i be as in the proof of Lemma 4.2. Construct the digraph F just as in the proof of Lemma 4.2, and remember that it consists of $m|X|$ pairwise edge disjoint copies of H which we called the essential copies of H in F (though it may well contain additional copies of H).

To construct D_1 which consists of digraphs that are ϵ -far from being H -free with high probability, we first construct F'_1 by removing each of the $m|X|$ essential copies of H , randomly and independently, with probability $1 - 1/|E(H)|$. We then create G_1 by taking an s blow up of F'_1 adding isolated vertices, if needed. Finally, D_1 consists of all randomly permuted copies of such digraphs G_1 . It follows from a standard Chernoff bound, that with probability $1 - o(1)$, at least $m|X|/2|E(H)|$ essential copies of H are left in F'_1 , where the $o(1)$ term tends to 0, as ϵ tends to 0. Similar to the derivation of (4), it is easy to show that if $m|X|/2|E(H)|$ of these copies of H are left in F'_1 , the graph G_1 is ϵ -far from being H free. It follows that with probability $1 - o(1)$, a member of D_1 is ϵ -far from being induced H -free.

The distribution D_2 of digraphs that are H -free, is defined by first constructing F'_2 by randomly, independently and uniformly picking from a each of the $m|X|$ essential copies of H a single edge, and removing all the other edges of that copy. We then create G_2 by taking an s blow up of F'_2 adding isolated vertices, if needed. Finally, D_2 consists of all randomly permuted copies of such digraphs G_2 . The main argument of Lemma 4.2, states that the graph F defined in the Lemma contains only triangles whose three edges belong to one of the copies essential copies of H . Hence, keeping a single edge from each of these copies results in a triangle free graph, and in particular all the graphs in G_2 are H -free.

Now consider a set of vertices S in G_1 (or G_2) and its natural projection to a subset of $V(F)$, which we also denote by S with a slight abuse of notation. Suppose S has the property that it does not contain more than two vertices from any one of the essential copies of H .

If this property holds, then each edge spanned by S is contained in a different essential copy of H . Therefore, each edge has probability $1/|E(H)|$ of being in F'_1 , and these probabilities are mutually independent. Similarly, each such edge has probability $1/|E(H)|$ of being in F'_2 and these probabilities are also mutually independent. It follows that sampling a digraph G from D_1 , and looking at the induced digraph on a set S with the above property, has *exactly* the same distribution as sampling a digraph G from D_2 , and looking at the induced digraph on S .

To complete the proof we have to show that no deterministic algorithm can distinguish between the distributions D_1 and D_2 with constant probability. To this end, it is clearly enough to show that with probability $1 - o(1)$, any deterministic algorithm that looks at a digraph spanned by less than $(1/\epsilon)^{c' \log 1/\epsilon}$ vertices, has *exactly* the same probability of seeing any digraph regardless of the distribution from which the digraph was chosen. By the discussion in the previous paragraph, this can be proved by establishing that, with high probability, a small set of vertices does not contain

three vertices from the same essential copy of H . For a fixed ordered set of three vertices in S , consider the event that they all belong to the same copy of H . The first two vertices determines all the vertices of one of these copies uniquely. Now, the conditional probability that the third vertex is also a vertex of the same copy is $h/|V(F)| \leq h/m$. By the union bound, the probability that the required property is violated, assuming $|S| = D$, is at most

$$hD^3/m \leq hD^3\epsilon^{c \log 1/\epsilon}.$$

This quantity is $o(1)$ as long as $D = o((1/\epsilon)^{\frac{c}{3} \log 1/\epsilon})$, where here we applied the lower bound on the size of m given in (3). Therefore, if the algorithm has query complexity $o((1/\epsilon)^{c' \log 1/\epsilon})$ for some absolute positive constant c' , it has probability $1 - o(1)$ of looking at a subset on which the distributions D_1 and D_2 are identical, thus, the probability that it distinguishes between D_1 and D_2 is $o(1)$. \blacksquare

A slightly more complicated argument than the above can give two distributions D_1 and D_2 , such that the graphs in D_1 are *always* ϵ -far from being induced H -free, while the graphs in D_2 are always H -free. The idea is to first partition the $m|X|$ essential copies of H into groups of size $|E(H)|$ assuming for simplicity that $|E(H)|$ divides $m|X|$. To create D_1 , we randomly pick from each group of $|E(H)|$ copies of H a single copy, and delete all its edges. To create D_2 , we do exactly the same as we did in the proof of Theorem 4. It is easy to appropriately modify the proof above in order to show that any deterministic algorithm with query complexity $o((1/\epsilon)^{c \log 1/\epsilon})$ can not distinguish between D_1 and D_2 . As this argument has no qualitative advantage, we described the simpler one given above.