

## ON SPARSE GRAPHS WITH DENSE LONG PATHS

P. ERDÖS\*, R. L. GRAHAM† and E. SZEMERÉDI\*  
Stanford University, Stanford, California, U.S.A.

(Received April 1975)

### INTRODUCTION

The following problem was raised by H.-J. Stoss [3] in connection with certain questions related to the complexity of Boolean functions. An acyclic directed graph  $G$  is said to have property  $P(m, n)$  if for any set  $X$  of  $m$  vertices of  $G$ , there is a directed path of length  $n$  in  $G$  which does not intersect  $X$ . Let  $f(m, n)$  denote the minimum number of edges a graph with property  $P(m, n)$  can have. The problem is to estimate  $f(m, n)$ .

In this paper we shall restrict ourselves to the case  $m = n$ . We shall prove

$$c_1 n \log n / \log \log n < f(n, n) < c_2 n \log n \quad (1)$$

(where  $c_1, c_2, \dots$ , will hereafter denote suitable positive constants). In fact, the graph we construct in order to establish the upper bound on  $f(n, n)$  in (1) will have just  $c_3 n$  vertices. In this case the upper bound in (1) is essentially best possible since it will also be shown that for  $c_4$  sufficiently large, if a graph on  $c_4 n$  vertices has property  $P(n, n)$  then it must have at least  $c_5 n \log n$  edges.

### A PRELIMINARY LEMMA

In order to establish the upper bound in (1) we first need the following result.

*Lemma.* For all  $\delta > 0$  there exists  $c = c(\delta)$  such that for all  $t$  sufficiently large, there exists a bipartite graph  $B = B(\delta; t)$  with vertex sets  $A$  and  $A'$  so that:

- (i)  $|A| = |A'| = t$ ;
- (ii)  $B$  has at most  $c(\delta)t$  edges;
- (iii) If  $X \subseteq A, X' \subseteq A'$  with  $|X| \geq \delta t, |X'| \geq \delta t$  then  $(X, X') = \{\{x, x'\} : x \in X, x' \in X'\}$  contains an edge of  $B$ .

*Proof:* We use a simple probabilistic argument to show the existence of  $B$ . Form a bipartite graph  $\bar{B}$  on the vertex sets  $A$  and  $A'$  with  $|A| = |A'| = t$  by selecting for each  $a \in A$  a random subset  $\bar{B}(a) \subseteq A'$  of cardinality  $d = d(\delta)$  (to be specified later). Call  $\bar{B}$  "bad" if there exists  $X \subseteq A, X' \subseteq A'$ , with  $|X| \geq \delta t, |X'| \geq \delta t$ , so that  $(X, X')$  contains no edge of  $\bar{B}$ . For fixed  $X$  and  $X'$ , the probability that  $\bar{B}$  is bad because of these two subsets is at most

$$\binom{(1-\delta)t}{d}^{d\delta} \binom{t}{d}^{(1-\delta)\delta} / \binom{t}{d}^d.$$

Hence, the total probability that  $\bar{B}$  is bad is at most

$$\binom{t}{\delta t}^2 \binom{(1-\delta)t}{d}^{d\delta} \binom{t}{d}^{(1-\delta)\delta} / \binom{t}{d}^d.$$

A simple computation shows that if  $d$  is chosen suitably large, for example, so that

$$(1 - \delta^2)^{d\delta} < 1/4,$$

\*P. Erdős and E. Szemerédi, Mathematical Institute of the Hungarian Academy of Sciences, Réaltanoda ut. 13-15 Budapest V, Hungary.

†R. L. Graham, Bell Laboratories, Department of Mathematics, 600 Mountain Avenue, Murray Hill, New Jersey, 07974, U.S.A.

This research was supported in part by National Science Foundation grant GJ 36473X, by the Office of Naval Research contract NR 044-402, and by IBM Corporation. Reproduction in whole or in part is permitted for any purpose of the United States Government.

then for  $t$  sufficiently large (e.g.,  $t > d/\delta$ ) this probability is less than 1, and so, a graph  $B = B(\delta; t)$  must exist which satisfies the requirements of the lemma.

CONSTRUCTION OF  $G$

The next step in the proof of (1) is the construction of the directed graph  $G$ . For large  $n$ ,  $G = G(n)$  will have as its vertex set  $V = \{0, 1, \dots, 2^n - 1\}$ . If  $v$  and  $m$  are positive integers, then  $D_v(m)$  will denote the set  $\{v, v + 1, \dots, v + m - 1\} \cap V$ . Similarly,  $D_v^*(m)$  will denote the set  $\{v, v - 1, \dots, v - m + 1\} \cap V$ . In general,  $\epsilon_1, \epsilon_2, \dots$ , will denote suitably chosen fixed positive constants to be specified later. The edge set  $E$  of  $G$  is formed as follows:

- (i) For  $v \in V$ , the pairs  $(v, x)$ ,  $x \in D_{v+1}(4n)$ , are in  $E$ ;
- (ii) For each  $t$  with  $n/2 \leq 2^t < 2^n$ , and each  $i$  as specified below a copy of  $B(\epsilon_1; 2^t)$  is formed between the vertex sets  $A = D_{m \cdot 2^t}(2^t)$  and  $A' = D_{(m+i) \cdot 2^t}(2^t)$ ,  $0 \leq m < 2^{n-t}$ , where  $i = 1, 2, \dots, 10$  (or if  $i$  cannot assume the value 10 because  $(m + 10)2^t > 2^n$ , then it ranges from 1 to  $2^{n-t} - m$ ). All edges are directed from  $x$  to  $y$  with  $x < y$ .

An elementary calculation shows that

$$|E| < c_6 n 2^n.$$

THE UPPER BOUND

*Theorem 1.* For a suitable  $\epsilon > 0$ ,  $G(n)$  has property  $P(\epsilon \cdot 2^n, \epsilon \cdot 2^n)$  for all sufficiently large  $n$ .

*Proof:* The theorem will be proved by a sequence of claims. First we show that  $G(n)$  shares with the graphs  $B(\epsilon; t)$  the following property.

*Claim 1.* If  $m \geq 2n$  and  $X \subseteq D_x(m)$ ,  $X' \subseteq D_{x+m}(m)$ , satisfy  $|X| \geq \epsilon_2 m$ ,  $|X'| \geq \epsilon_2 m$ , then  $[X, X'] = \{(x, x') : x \in X, x' \in X'\}$  contains an edge of  $G(n)$ .

*Proof of Claim:* Let  $2^t \leq m/2 < 2^{t+1}$ . Thus,  $m/4 < 2^t$  so at most five of the intervals  $D_{r \cdot 2^t}(2^t)$  intersect  $D_x(m)$  and at most five of them intersect  $D_{x+m}(m)$ . Since  $|X| \geq \epsilon_2 m$  then some  $D_{r \cdot 2^t}(2^t)$  and  $D_{r' \cdot 2^t}(2^t)$  have

$$|D_{r \cdot 2^t}(2^t) \cap X| \geq \epsilon_2 m/5, |D_{r' \cdot 2^t}(2^t) \cap X'| \geq \epsilon_2 m/5. \tag{3}$$

But we must have  $|r' - r| \leq 10$  so that by the construction of  $G(n)$  there is a copy of  $B(\epsilon_1; 2^t)$  between  $D_{r \cdot 2^t}(2^t)$  and  $D_{r' \cdot 2^t}(2^t)$ . Thus, if  $\epsilon_2/5 > \epsilon_1$  and  $m \geq 2^t$  then the property of  $B(\epsilon_1; 2^t)$  guaranteed by the Lemma implies that  $[X, X']$  contains an edge of  $G(n)$  provided that  $t$  is sufficiently large (which is guaranteed by choosing  $n$  large enough). This proves the claim.

Next, let us choose an arbitrary fixed set  $X$  of vertices with  $|X| \leq \epsilon \cdot 2^n$ . The vertices in  $X$  will be referred to as the *marked* vertices of  $G$ ; the remaining vertices of  $G$  will be called the *unmarked* vertices of  $G$ .

Let us call an unmarked vertex  $y \in V$  *bad* if for some  $m \geq 1$  either at least  $\epsilon_3 m$  vertices in  $D_y(m)$  are marked or at least  $\epsilon_3 m$  vertices in  $D_y^*(m)$  are marked. Otherwise, an unmarked vertex of  $G$  is called *good*.

*Claim 2.* There are at most  $\epsilon_4 2^n$  bad vertices.

*Proof of Claim:* Let  $y_1$  denote the least unmarked vertex of  $G$  (if it exists) for which for some  $m_1 \geq 1$ , at least  $\epsilon_3 m_1$  vertices in  $D_{y_1}(m_1)$  are marked. In general, if  $y_1, \dots, y_k$  and  $m_1, \dots, m_k$  have been defined, let  $y_{k+1}$  be the least unmarked vertex of  $G$  following  $y_k + m_k - 1$  (if it exists) for which for some  $m_{k+1} \geq 1$  at least  $\epsilon_3 m_{k+1}$  vertices in  $D_{y_{k+1}}(m_{k+1})$  are marked. We continue this process until it no longer can be applied, so that, say,  $y_1, \dots, y_s$  and  $m_1, \dots, m_s$  have been defined. Similarly, let  $y_k^*$  denote the greatest unmarked vertex (if it exists) for which for some  $m_k^* \geq 1$ , at least  $\epsilon_3 m_k^*$  vertices in  $D_{y_k^*}^*(m_k^*)$  are marked, etc. In this way, we define  $y_1^*, \dots, y_s^*$  and  $m_1^*, \dots, m_s^*$ .

It follows from the preceding construction and the definition of a bad vertex that *all* bad vertices are contained in the set

$$Y = \bigcup_{k=1}^s D_{y_k}(m_k) \cup \bigcup_{k=1}^{s^*} D_{y_k^*}^*(m_k^*)$$

Thus, there are at most

$$M = \sum_{k=1}^s m_k + \sum_{k=1}^{s^*} m_k^*$$

bad vertices. However, by our construction there are at least  $(\epsilon_3/2)M$  marked vertices in  $Y$ . Since by hypothesis there are at most  $\epsilon \cdot 2^n$  marked vertices in  $V$  then we have

$$\begin{aligned} (\epsilon_3/2)M &\leq \epsilon \cdot 2^n, \\ M &\leq (2\epsilon/\epsilon_3)2^n < \epsilon_4 2^n, \end{aligned}$$

which proves the claim.

For an unmarked vertex  $x$ , let  $P_x(m)$  denote the set of all unmarked vertices in  $D_x(m)$  which can be reached from  $x$  by directed paths which contain only unmarked vertices.

*Claim 3.* If  $x$  is a good vertex and  $|D_x(m)| = m$  then

$$|P_x(m)| > \epsilon_5 m \tag{4}$$

*Proof of Claim:* If  $m \leq 4n$  then since  $x$  is good, at least  $(1 - \epsilon_3)m$  vertices in  $D_x(m)$  are unmarked and  $x$  has edges directly to all of them. Suppose  $m > 4n$ . Let  $m'$  denote  $\lceil m/2 \rceil$ . Since  $|D_x(m')| = m'$  then by induction  $|P_x(m')| > \epsilon_5 m'$ . Since  $x$  is good then at most  $\epsilon_3 m$  vertices in  $D_x(m)$  are marked. Hence, at most  $\epsilon_3 m$  vertices in  $D_{x+m'}(m') \subseteq D_x(m)$  are marked. Since  $m' \geq 2n$  and  $\epsilon_5 \geq \epsilon_2$  then there are edges from  $P_x(m')$  to at least  $(1 - \epsilon_2)m'$  vertices in  $D_{x+m'}(m')$ . But at most  $\epsilon_3 m < 3\epsilon_3 m'$  vertices in  $D_{x+m'}(m')$  are marked. Hence,  $P_x(m')$  must have edges to at least  $(1 - \epsilon_2 - 3\epsilon_3)m'$  unmarked vertices in  $D_{x+m'}(m')$ . Since  $1 - \epsilon_2 - 3\epsilon_3 > 3\epsilon_5$  then

$$|P_x(m)| > 3\epsilon_5 m' > \epsilon_5 m.$$

The claim now follows by induction.

In exactly the same way it follows that if  $P_x^*(m)$  denotes the set of all unmarked vertices in  $D_x^*(m)$  which are connected to the unmarked vertex  $x$  by a directed path containing only unmarked vertices, and  $x$  is a good vertex and  $|D_x^*(m)| = m$ , then

$$|P_x^*(m)| > \epsilon_5 m. \tag{4}$$

*Claim 4.* Let  $x$  and  $x'$  be good vertices with  $x < x'$ . Then  $x' \in P_x(2^n)$ .

*Proof:* If  $x' - x \leq 4n$  then the claim is immediate since by construction there is an edge from  $x$  to  $x'$ . Assume  $x' - x > 4n$ . Let  $y = \lceil (x + x')/2 \rceil$  and let  $m = y - x + 1$ . Consider the intervals  $D_x(m)$  and  $D_{x'}^*(m)$ . Either they are adjacent or they have the single element  $y$  in common. Since  $x$  and  $x'$  are good then by (4) and (4')

$$|P_x(m)| > \epsilon_5 m, |P_{x'}^*(m)| > \epsilon_5 m. \tag{5}$$

Since  $\epsilon_5 \geq \epsilon_2$  then by Claim 1, there is an edge in  $G$  from a vertex of  $P_x(m)$  to a vertex of  $P_{x'}^*(m)$ . Thus, there is a directed path from  $x$  to  $x'$  containing no marked vertices and the claim is proved.

The proof of the theorem is now immediate. By Claim 2 there are at least  $(1 - \epsilon_4 - \epsilon)2^n$  good vertices in  $G$ . By Claim 4 we can form a directed path which contains only unmarked vertices and which contains *all* the good vertices (since  $x'$  can always be chosen to be the next good vertex following  $x$ ). Since  $1 - \epsilon_4 - \epsilon > \epsilon$  then the theorem follows (where it is easily seen how the appropriate values of  $\epsilon_k$  and  $c_k$  can be chosen).

#### THE LOWER BOUND

The following result will establish the lower bound in (1).

*Theorem 2.* Let  $H$  be an acyclic directed graph with at most  $c_7 n \log n / \log \log n$  edges where  $n$  is a large fixed integer. Then there is a set of at most  $n$  vertices of  $H$  which hits every directed path of length  $n$ .

*Proof:* Let us denote the vertex set of  $H$  by  $V = \{1, 2, \dots, v\}$ . We may assume that all edges are of the form  $(i, j)$  with  $i < j$ . For an edge  $e = (i, j)$  of  $H$ , let  $length(e)$  be defined to be  $j - i$ .

Partition the edges of  $H$  into classes  $C_0, C_1, \dots, C_r$  where

$$C_k = \{e: 2^{4k \log \log n} \leq \text{length}(e) < 2^{4(k+1) \log \log n}\}$$

and  $r = \lceil \log v/4 \log \log n \rceil$ .

Since  $H$  has at least  $c_8 n \log n / \log \log n$  edges then it follows that  $v \geq c_9 n^{1/2}$  and  $r \geq c_{10} \log n / \log \log n$ . Hence some class  $C_a$  with  $0 \leq a < r$  has at most  $c_{11} n$  elements. Let us delete all vertices in  $H$  incident to any of the edges in  $C_a$ . Furthermore, we also delete those vertices  $x \in V$  which satisfy

$$0 \leq x - m \cdot 2^{4a \log \log n} (1 + 2^{2 \log \log n}) < 2^{4a \log \log n}$$

for some integer  $m \geq 0$ . This latter step removes at most

$$\left(\frac{2}{2^{2 \log \log n} - 1}\right)v = O(n)$$

vertices, since  $v \leq 2c_7 n \log n / \log \log n$ . Hence we have deleted at most  $c_{12} n$  vertices altogether. However, any directed path remaining has at most

$$\left(\frac{2^{(4a+2) \log \log n} - 2^{4a \log \log n}}{2^{4(a+1) \log \log n}}\right)v = O(n)$$

edges, since we cannot go more than  $2^{(4a+2) \log \log n} - 2^{4a \log \log n}$  steps without using an edge whose length exceeds  $2^{4a \log \log n}$ ; and the length of such an edge actually exceeds  $2^{4(a+1) \log \log n}$ . This proves the theorem.

By using a different partition of the edges of  $H$ , namely, into the classes  $C'_0, \dots, C'_r$  where

$$C'_k = \{e: 2^{c_{13} k} \leq \text{length}(e) < 2^{c_{13}(k+1)}\}$$

for a suitable constant  $c_{13}$ , we can establish the following result.

*Theorem 3.* If  $c_{14}$  is sufficiently large then any graph  $G$  on  $c_{14} n$  vertices having property  $P(n, n)$  must have at least  $c_{15} n \log n$  edges.

The graphs  $G(n)$  used in Theorem 1 show that the result in Theorem 3 is to within constant factors best possible.

#### SOME RELATED QUESTIONS

We now consider several problems for ordinary (undirected) graphs. Let  $F_e(n, n)$  (resp.,  $F_v(n, n)$ ) denote the smallest integer for which there is a graph with  $F_e(n, n)$  edges so that the deletion of any  $n$  of its vertices there still remains a connected component of  $n$  edges (resp., vertices). We shall prove by probabilistic methods that

$$F_e(n, n) < c_{16} n, F_v(n, n) < c_{17} n. \tag{6}$$

The method we use is the same as that in the work of Erdős and Renyi [1], [2]. It turns out that almost all graphs have the desired property.

*Theorem 4.* For every  $\epsilon > 0$  there is a  $c = c(\epsilon)$  so that all but  $O\left(\binom{\binom{2+\epsilon}{cn}}{cn}\right)$  graphs  $G$  with  $(2 + \epsilon)n$  vertices and  $cn$  edges have the property that after the omission of any  $n$  of its vertices, a connected component of at least  $n$  vertices remains.

*Proof:* It suffices to show that if  $n$  vertices are omitted and the remaining  $n(1 + \epsilon)$  vertices are split into two classes  $S_1$  and  $S_2$  with  $|S_1| \geq \epsilon n$ ,  $|S_2| \geq \epsilon n$ , then there is at least one edge joining a vertex of  $S_1$  to a vertex of  $S_2$ .

Consider a random graph  $G$  on  $(2 + \epsilon)n$  vertices and  $cn$  edges (where  $c$  will be specified later). There are  $\binom{(2 + \epsilon)n}{n}$  ways that  $n$  vertices of  $G$  can be deleted. The remaining  $n(1 + \epsilon)$  points

can then be split into two sets  $S_1$  and  $S_2$  in at most  $2^{n(1+\epsilon)}$  ways. Thus, the total number of splittings is at most

$$\binom{(2+\epsilon)n}{n} 2^{n(1+\epsilon)} < 2^{(2+\epsilon)n} 2^{n(1+\epsilon)} < 2^{3(1+\epsilon)n}$$

Between  $S_1$  and  $S_2$  there are at least  $\epsilon n^2$  potential edges. The probability that none of these edges actually occurs in  $G$  is less than  $\left(1 - \frac{c}{(2+\epsilon)n}\right)^{\epsilon n^2}$ . Thus, if  $c$  is chosen so that

$$2^{3(1+\epsilon)n} \left(1 - \frac{c}{(2+\epsilon)n}\right)^{\epsilon n^2} \rightarrow 0 \tag{7}$$

as  $n \rightarrow \infty$  then we easily see that almost all graphs cannot be split in such a way.

Since

$$\left(1 - \frac{c}{(2+\epsilon)n}\right)^{\epsilon n^2} \rightarrow e^{-(\epsilon c/(2+\epsilon))n}$$

then for  $c$  large enough, e.g.,  $c > 18(\epsilon + \epsilon^{-1})$ ,

$$e^{-(\epsilon c/(2+\epsilon))n} < e^{-3(1+\epsilon)n}$$

and (7) holds. This proves the theorem.

The other half of (6) is proved in a similar way. It would be interesting to determine the best possible value of  $c$  but this does seem to be too easy.

We mention here the undirected analogue of (1). Let  $g(n, n)$  denote the smallest integer for which there is an undirected graph of  $g(n, n)$  edges so that if we omit any  $n$  of its vertices then there always remains a path of length  $n$ . We believe

$$\frac{g(n, n)}{n} \rightarrow \infty, \quad \frac{g(n, n)}{n \log n} \rightarrow 0$$

as  $n \rightarrow \infty$  and hope to return to this question in finite time.

A related question is the following: Consider random graphs on  $n$  vertices and  $Cn$  edges. Is it true that for large  $C$  almost all of these graphs have a path of length  $n(1 - \epsilon)$ ? It is known[4] that almost all graphs on  $n$  vertices and  $(1/2 + \epsilon)n \log n$  edges are Hamiltonian.

It is possible to introduce another parameter into these questions. Let  $F_v(t; n, n)$  denote the smallest integer for which there is a graph with  $t$  vertices and  $F_v(t; n, n)$  edges having the property that if any  $n$  vertices are deleted there still remains a connected component with at least  $n$  vertices. If  $t/n \rightarrow c > 2$  then  $F_v(t; n, n)/n \rightarrow A(c)$  where  $A(c) \rightarrow \infty$  as  $c \rightarrow 2$ . (The behavior of  $F_e(t; n, n)/n$  is similar). We would also omit edges instead of vertices but leave the formulation of these questions to the reader.

*Acknowledgment*—The authors gratefully acknowledge the useful suggestions on this problems given to us by D. E. Knuth.

REFERENCES

1. P. Erdős and A. Rényi, On random graphs I, *Publ. math. Debrecen* 6 290–297, (1959).
2. P. Erdős and A. Rényi, On the evolution of random graphs, *Bull. int. stat. Inst.* 38 343–347 (1961).
3. D. E. Knuth, (personal communication).
4. J. Komlós and E. Szemerédi, On Hamiltonian Circuits in Random Graphs, (to appear).