

for such i , $\binom{i}{j+2^{n-1}} = \binom{i}{j}$. By Lemma 3, the well-known combinatorial identity

$$\binom{i+2^{n-1}}{j+2^{n-1}} = \sum_{k=0}^{j+2^{n-1}} \binom{i}{k} \binom{2^{n-1}}{j+2^{n-1}-k}$$

reduces in GF(2) to

$$\binom{i+2^{n-1}}{j+2^{n-1}} = \binom{i}{j} + \binom{i}{j+2^{n-1}}. \tag{1}$$

But $\binom{i}{j+2^{n-1}}$ is periodic with period 2^n , so that

$$\binom{i+2^{n-1}}{j+2^{n-1}} = \binom{i-2^{n-1}}{j+2^{n-1}}. \tag{2}$$

Note that $0 \leq i - 2^{n-1} < 2^{n-1} \leq j + 2^{n-1}$, so that the right-hand side of (2) is 0, and (1) becomes, in GF(2),

$$\binom{i}{j+2^{n-1}} = \binom{i}{j}. \quad \square$$

Corollary 5: For $\mathbf{d} = (d_0, d_1, \dots, d_{2^{n-1}-1}) \neq \mathbf{0}$, let $\mathbf{d}' = (\mathbf{0} : \mathbf{d})$, $\mathbf{0} \in \text{GF}(2)^{2^{n-1}}$, then $(\mathbf{d}') \in B(2^n)$ has $c(\mathbf{d}') = 2^{n-1} + c(\mathbf{d})$.

Proof: For $i \geq 0$, suppose

$$d_i = \sum_{j=0}^{c(\mathbf{d})-1} a_j \binom{i}{j}$$

with $1 \leq c(\mathbf{d}) \leq 2^{n-1}$ and $a_{c(\mathbf{d})-1} = 1$. Then Lemma 4 implies

$$d'_i = \sum_{j=0}^{c(\mathbf{d})-1} a_j \binom{i}{j+2^{n-1}}. \quad \square$$

Theorem 6: Let $\mathbf{s} = (s_0, s_1, \dots, s_{2^n-1}) = (L(\mathbf{s}) : R(\mathbf{s}))$ be a binary vector with associated sequence $(s) \in B(2^n)$, $n \geq 1$. Form $\mathbf{d} = R(\mathbf{s}) - L(\mathbf{s})$, so that $(\mathbf{d}) \in B(2^{n-1})$.

- a) If $\mathbf{d} = \mathbf{0}$, then $c(\mathbf{s}) = c(L(\mathbf{s}))$.
- b) If $\mathbf{d} \neq \mathbf{0}$, then $c(\mathbf{s}) = 2^{n-1} + c(\mathbf{d})$.

Proof: If $\mathbf{d} = \mathbf{0}$, then $R(\mathbf{s}) = L(\mathbf{s})$, so $(\mathbf{s}) = ((L(\mathbf{s}) : L(\mathbf{s}))) = (L(\mathbf{s}))$ and $c(\mathbf{s}) = c(L(\mathbf{s}))$. If $\mathbf{d} \neq \mathbf{0}$, then write

$$\begin{aligned} \mathbf{s} &= (L(\mathbf{s}) : R(\mathbf{s})) \\ &= (L(\mathbf{s}) : L(\mathbf{s})) + (\mathbf{0} : R(\mathbf{s}) - L(\mathbf{s})) \\ &= (L(\mathbf{s}) : L(\mathbf{s})) + (\mathbf{0} : \mathbf{d}). \end{aligned}$$

Let $\mathbf{d}' = (\mathbf{0} : \mathbf{d})$, so that $(\mathbf{s}) = (L(\mathbf{s})) + (\mathbf{d}')$. Note that $c(L(\mathbf{s})) \leq 2^{n-1}$ and, by Corollary 5, $c(\mathbf{d}') = 2^{n-1} + c(\mathbf{d}) > 2^{n-1}$, so that the term $\binom{i}{c(\mathbf{d}')-1+2^{n-1}}$ occurs in the basis expansion of $s_i = L(\mathbf{s})_i + d'_i$. Thus $c(\mathbf{s}) = c(\mathbf{d}') = 2^{n-1} + c(\mathbf{d})$. \square

The algorithm of Section II applies the result of Theorem 6 recursively, starting with the initial vector \mathbf{s} of length 2^n , and stopping when the vector $(\mathbf{0})$ (of complexity 0) or the vector (1) (of complexity 1) is encountered.

CONCLUSION

For a given sequence of arbitrary period N , the Massey algorithm [3] accepts the sequence sequentially and at each stage computes the connection polynomial for the shortest LFSR that generates the encountered portion of the sequence. The Massey

algorithm may have to run through more than one period of length N of the sequence before it stabilizes on the correct connection polynomial. In practice, additional iterations are required to ensure that the algorithm has in fact stabilized. The algorithm given in this correspondence works only for a sequence with period of length $N = 2^n$ and computes the complexity c in $\log N = n$ steps. The connection polynomial $f(E)$ then must be $(E - 1)^c$ in this case. The storage requirements of the Massey algorithm depend directly on the eventual complexity of the sequence, while the present algorithm must always store a single period of the sequence, making the algorithm inappropriate for very long periods.

REFERENCES

- [1] E. J. Groth, "Generation of binary sequences with controllable complexity," *IEEE Trans. Inform. Theory*, vol. IT-17, pp. 288-296, 1971.
- [2] E. L. Key, "An analysis of the structure and complexity of non-linear binary sequence generators," *IEEE Trans. Inform. Theory*, vol. IT-22, pp. 732-736, Nov. 1976.
- [3] J. L. Massey, "Shift registers synthesis and BCH decoding," *IEEE Trans. Inform. Theory*, vol. IT-15, pp. 122-127, Jan. 1969.
- [4] S. W. Golomb, *Shift Register Sequences*. San Francisco: Holden-Day, 1967.

The Largest Super-Increasing Subset of a Random Set

EHUD D. KARNIN, STUDENT MEMBER, IEEE, AND
MARTIN E. HELLMAN, FELLOW, IEEE

Abstract—It is shown that the longest super-increasing sequence which can be constructed from a set of n independent uniformly distributed random variables is almost surely asymptotic to $\log_2 n$. Some extensions of this result, as well as the implications for the security of knapsack-based cryptographic systems, are discussed.

I. INTRODUCTION AND MOTIVATION

Given a set A of n positive numbers u_i , $i = 1, 2, \dots, n$ and a positive number S , the knapsack problem [1] is to find a binary solution vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$ such that

$$\sum_{i=1}^n x_i u_i = S \quad x_i \in \{0, 1\}. \tag{1}$$

The associated decision problem, i.e., determining whether (1) has a solution, is in the class of NP-complete problems [2]. Therefore, it is believed that solving (1) is very hard in general. However, there are sets A for which the problem is extremely easy.

Consider the case when the sequence u_i , $i = 1, 2, \dots, n$ is *super-increasing*:

$$u_j > \sum_{i=1}^{j-1} u_i \quad j = 2, 3, \dots, n. \tag{2}$$

When (2) holds, (1) is solvable in time which grows only linearly in n (cf. [1]). Also when a large subset of A can be ordered as a super-increasing sequence, the effort of solving (1) is substantially reduced (as discussed in Section V).

The knapsack problem forms the basis for a public key cryptographic system [1]. The public key is generated by choosing a

Manuscript received December 7, 1981; revised March 7, 1982. This research was supported in part by the National Security Agency under Contract MDA904-81-C-0414, and in part by Joint Services Electronics Program under Contract DAAG29-81-0057.

The authors are with the Department of Electrical Engineering, Informations System Laboratory, Stanford University, Stanford, CA 94305.

sequence of n super-increasing integers, multiplying each by w and reducing modulo m (where w and m are large integers, typically of the order of 2^{2^n}). This transformation is the basis of most pseudorandom number generators and produces integers which appear to be uniformly distributed in $\{0, 1, \dots, m\}$. Dividing by m and neglecting quantization effects (since m is very large), we can model the resulting u_i , $i = 1, 2, \dots, n$ as being uniformly distributed in $[0, 1]$.

Motivated by this we consider the following question: let $A = \{u_i\}_{i=1}^n$ with each u_i independent and uniformly distributed in $[0, 1]$. Look at all subsets of A which have elements that can be ordered as a super-increasing sequence. What is the cardinality μ_n of the largest such super-increasing subset?

This problem and some extensions are studied in the following sections.

II. MAIN RESULT

Theorem:

$$\frac{\mu_n}{\log_2 n} \rightarrow 1 \quad \text{a.s.} \quad (3)$$

Proof: We establish a lower and an upper bound for $\mu_n/\log_2 n$.

Lower Bound: Consider the events

$$A_i = \{\exists u_i \in [2^{-(i-1)} - 2^{-m}, 2^{-(i-1)}]\}, \quad i = 1, 2, \dots, m.$$

Then

$$P(\mu_n \geq m) \geq P\left(\bigcap_{i=1}^m A_i\right),$$

which is equivalent to

$$P(\mu_n < m) \leq P\left(\bigcup_{i=1}^m A_i^c\right) \leq \sum_{i=1}^m P(A_i^c) = m(1 - 2^{-m})^n.$$

Let $m = (1 - \epsilon)\log_2 n$, $\epsilon > 0$. Then

$$\begin{aligned} P\left(\frac{\mu_n}{\log_2 n} < 1 - \epsilon\right) &\leq [(1 - \epsilon)\log_2 n][1 - 2^{-(1-\epsilon)\log_2 n}]^n \\ &= [(1 - \epsilon)\log_2 n] \left[\left(1 - \frac{1}{n^{1-\epsilon}}\right)^{n^{1-\epsilon}} \right] \\ &\rightarrow [(1 - \epsilon)\log_2 n] e^{-n^\epsilon} \rightarrow 0. \end{aligned}$$

Moreover these terms are summable, i.e.,

$$\sum_{n=2}^{\infty} P\left(\frac{\mu_n}{\log_2 n} < 1 - \epsilon\right) < \infty.$$

Consequently, by the Borel-Cantelli lemma

$$P\left(\liminf \frac{\mu_n}{\log_2 n} < 1 - \epsilon\right) = 0.$$

Since ϵ is arbitrary we get

$$\liminf \frac{\mu_n}{\log_2 n} \geq 1 \quad \text{a.s.} \quad (4)$$

Upper bound:

$$P(\mu_n \geq m) \leq P\{\exists u_i \text{ such that } u_i \leq 2^{-m}\}. \quad (5)$$

We prove this claim by contradiction. From (2), any super-increasing sequence starting with u_i should grow faster than

$$u_i, u_i, 2u_i, 4u_i, \dots, 2^{m-2}u_i.$$

But $u_i > 2^{2-m}$ implies that the last term is greater than one, a

contradiction. Relation (5) can be further bounded by

$$P(\mu_n \geq m) \leq n2^{2-m} = 4n2^{-m}.$$

Let $m = (1 + \epsilon)\log_2 n$, $\epsilon > 0$. Then

$$P\left(\frac{\mu_n}{\log_2 n} \geq 1 + \epsilon\right) \leq 4n2^{-(1+\epsilon)\log_2 n} = 4n^{-\epsilon} \rightarrow 0.$$

To establish the almost surely convergence consider the subsequence $n_i = 2^i$, and observe that

$$\sum_{i=1}^{\infty} P\left(\frac{\mu_{n_i}}{\log_2 n_i} \geq 1 + \epsilon\right) \leq \sum_{i=1}^{\infty} 4 \cdot 2^{-\epsilon i} < \infty.$$

Hence by the Borel-Cantelli lemma,

$$P\left(\limsup_{i \rightarrow \infty} \frac{\mu_{n_i}}{\log_2 n_i} \geq 1 + \epsilon\right) = 0$$

which implies, since ϵ is arbitrary, that

$$\limsup_{i \rightarrow \infty} \frac{\mu_{n_i}}{\log_2 n_i} \leq 1 \quad \text{a.s.}$$

Now consider n in the gap $2^{i-1} < n \leq 2^i$. The cardinality μ_n is monotonically increasing, so $\mu_n \leq \mu_{2^i}$. Also, $\log_2 n > \log_2 2^{i-1} = i - 1$. Hence,

$$\begin{aligned} \frac{\mu_n}{\log_2 n} &\leq \frac{\mu_{2^i}}{i-1} = \frac{\mu_{2^i}}{i(1-1/i)} \\ &\leq \frac{\mu_{2^i}}{\log_2 2^i} (1-1/i_0)^{-1}, \quad i \geq i_0. \end{aligned}$$

Therefore

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{\mu_n}{\log_2 n} &\leq \limsup_{i \rightarrow \infty} \frac{\mu_{n_i}}{\log_2 n_i} (1-1/i_0)^{-1} \\ &\stackrel{\text{a.s.}}{\leq} (1-1/i_0)^{-1}. \end{aligned}$$

Since i_0 may be taken arbitrarily large, this implies

$$\limsup \frac{\mu_n}{\log_2 n} \leq 1 \quad \text{a.s.} \quad (6)$$

Finally, (4) and (6) yield (3), which completes the proof.

III. HIGHER DIMENSIONS

In a k -dimensional vector space we define a relation: $\mathbf{a} > \mathbf{b}$ if each component of \mathbf{a} is greater than the corresponding component of \mathbf{b} . This relation is a partial ordering, and we use it to define a super-increasing sequence of vectors as

$$\mathbf{u}_j > \sum_{i=1}^{j-1} \mathbf{u}_i, \quad j = 2, 3, \dots, n, \quad (7)$$

where addition of vectors is, as usual, on a component by component basis.

As before let $\{u_i\}_{i=1}^n$ be independent, uniformly distributed in $[0, 1]^k$, and let μ_n be the length of the largest super-increasing sequence which can be constructed from them.

Corollary 1:

$$\frac{\mu_n}{\log_2 n} \rightarrow \frac{1}{k} \quad \text{a.s.} \quad (8)$$

Proof: For a lower bound take

$$A_i = \{\exists \mathbf{u}_i \in [2^{-(i-1)} - 2^{-m}, 2^{-(i-1)}]^k\}$$

and proceed as in the proof of the theorem. For the upper bound,

observe that

$$P(\mu_n \geq m) \leq P(\exists u_i \text{ such that } u_i < (2^{2-m}, 2^{2-m}, \dots, 2^{2-m}))$$

and then proceed as before.

IV. NONBINARY KNAPSACK

Suppose we modify the knapsack problem by asking for a solution of (1) such that

$$x_i \in \{0, 1, \dots, B-1\}$$

instead of $x_i \in \{0, 1\}$. Such "compact" knapsacks are of interest in cryptography, because they allow a reduction in the size of the public key [1]. The analog of an easily solved super-increasing knapsack is now

$$u_j > (B-1) \sum_{i=1}^{j-1} u_i \quad j = 2, 3, \dots, n, \quad (9)$$

which clearly reduces to (2) for the special case $B = 2$. When (9) holds the non binary knapsack is easily solved, so we ask about μ_n , the cardinality of the largest subset satisfying (9). Again we assume u_i independent, uniformly distributed in $[0, 1]$, and obtain the following.

Corollary 2:

$$\frac{\mu_n}{\log_B n} \rightarrow 1 \quad \text{a.s.} \quad (10)$$

Proof: For a lower bound consider the events

$$A_i = \{\exists u_i \in [B^{-(i-1)} - B^{-m}, B^{-(i-1)}]\}$$

and proceed as in the theorem. For the upper bound, observe that

$$P(\mu_n \geq m) \leq P(\exists u_i \text{ such that } u_i \leq \frac{B^2}{B-1} B^{-m})$$

and proceed as before. (For the subsequence argument take $n_i = B^i$.)

Note: The extension of this result to a k -dimensional vector space (cf. Section III) is

$$\frac{\mu_n}{\log_B n} \rightarrow \frac{1}{k} \quad \text{a.s.} \quad (11)$$

V. APPLICATION AND DISCUSSION

In cryptography a knapsack can be used to conceal an n bit message x (see [1]). The cryptanalyst knows u_i , $i = 1, 2, \dots, n$ and S , and tries to find x . Searching over all the 2^n possible x will yield a solution for (1). However, one may also use the following procedure.

For the $n - \mu_n$ u_i 's, which are not in the largest super-increasing subset, try all the $2^{n-\mu_n}$ possible x_i . In each trial subtract from S the u_i which correspond to $x_i = 1$. Then see if there is a solution to the associated easy super-increasing knapsack with μ_n components, in $O(\mu_n)$ operations.

By our theorem μ_n is almost surely asymptotic to $\log_2 n$, and the effort of solving (1) is reduced only by a factor of about $2^{\mu_n} = n$, compared to an exhaustive search. Thus the security of a knapsack-based cryptographic system is not substantially decreased by the above procedure.

As a final remark we compare a binary knapsack of length n and a nonbinary knapsack of length N . If both have to conceal the same amount of information, then $n = Nb$, where $b = \log_2 B$.

About $\log_2 n$ bits of information are involved in an "easy" super-increasing subset of the binary knapsack. For the nonbi-

nary knapsack we have by Corollary 2

$$\mu_N \sim \log_B N = \frac{\log_2(n/b)}{\log_2 B} \sim \frac{\log_2 n}{b}$$

Since each element carries b bits of information, $b\mu_N \sim \log_2 n$ bits of information are involved in the super-increasing part, which is the same number obtained for the binary knapsack.

ACKNOWLEDGMENT

The authors are thankful to Professor J. Michael Steele of Stanford University, Department of Statistics, for helpful discussions.

REFERENCES

- [1] R. Merkle and M. Hellman, "Hiding information and signatures in trapdoor knapsacks," *IEEE Trans. Inform. Theory*, vol. IT-24, pp. 525-530, Sept. 1978.
- [2] M. Garey and D. Johnson, *Computers and Intractability—A Guide to the Theory of NP-Completeness*. San Francisco: W. H. Freeman, 1979.

A Combinatorial Approach to Polygon Similarity

DAVID AVIS AND HOSSAM ELGINDY

Abstract—A new approach is presented to the classification problem of planar shapes represented by polygons. A shape is abstracted combinatorially by means of its visibility graph, and two shapes are deemed similar whenever their graphs are cyclically isomorphic. Efficient algorithms are presented for performing these operations.

I. INTRODUCTION

In pattern recognition, a common and economical approach to representing planar shapes within a computer is to describe their boundaries by polygons with a finite number of vertices [1, pp. 168-184]. With appropriate selection of vertices on or near the boundary, the shape can be adequately described without affecting its important features. Such polygons may be described by a sequence of edges representing a piecewise approximation to the boundary, or a graph which preserves important properties of the original shape.

The problem of classification requires the partition of the set of all polygons into a number of classes such that the polygons in each class are equivalent (or similar) with respect to a selected equivalence relation. In other words, a polygon can be uniformly mapped onto other polygons in its equivalence class such that certain features are preserved.

A hierarchical approach to the classification problem would involve classifying the polygons into broad equivalence classes, and then using further discrimination criteria within each class.

In this correspondence, we present a purely combinatorial technique. We define a finite graph on the vertices of the polygon and then deal exclusively with the graph. Two polygons are said to be equivalent whenever the corresponding graphs are isomorphic under some cyclic permutation of the vertices.

This equivalence relation decomposes the class of n -vertex polygons into broad equivalence classes. For example, it groups all n -vertex convex polygons in one class, which may be thought of as the structurally most simple class. Further, and more

Manuscript received May 14, 1981; revised April 20, 1982. This work was supported in part by Natural Sciences and Engineering Research Council (N.S.E.R.C.) Grant A3013 and in part by Fonds de Chercheurs et Action Concertée (F.C.A.C.) Grant EQ-1678.

The authors are with the School of Computer Science, McGill University, 805 Sherbrooke Street West, Montreal, PQ, Canada, H3A 2K6.