

Kolmogorov complexity of enumerating finite sets

Nikolai K. Vereshchagin*

Keywords: Kolmogorov complexity, the a priori probability.

Abstract

Solovay [5] has proved that the minimal length of a program enumerating a set A is upper bounded by 3 times the negative logarithm of the probability that a random program will enumerate A . It is unknown whether one can replace the constant 3 by a smaller constant. In this paper, we show that the constant 3 can be replaced by the constant 2 for *finite* sets A .

We recall first two complexity measures (“information content”) of computably enumerable sets attributed by Solovay in [5] to G. Chaitin (wee keep Solovay’s notations).

Let M be a machine with one infinite input tape and one infinite output tape. At the start the input tape contains an infinite binary string ω called the input to M . The output tape is empty at the start. We say that a program p *enumerates* a set $A \subset \mathbb{N} = \{1, 2, \dots\}$ if in the run on every input ω extending p machine M prints all the elements of A in some order and no other elements, and does not move the head on input tape beyond p . We do not require M to halt in the case when A is finite.¹ Let $I_M(A)$ denote the minimal length of a program enumerating A . There is a machine M_0 (called a *universal* machine) such that for every other machine M there is a constant c such that

$$I_{M_0}(A) \leq I_M(A) + c$$

for all $A \subset \mathbb{N}$. Fix any such M_0 and call $I(A) \stackrel{def}{=} I_{M_0}(A)$ the *complexity of enumeration of A* . This complexity thus depends on the choice of the universal machine but this dependence is rather weak: for any other universal machine M_1 the difference $|I_{M_0}(A) - I_{M_1}(A)|$ is bounded by a constant not depending on A .

The second complexity measure is related to the a priori probability distribution on enumerable sets. The definitions are as follows. Let M be a machine

*Moscow State University, Leninskie Gory 1, Moscow 119992, Email: ver@mccme.ru. Work was done while visiting Laboratoire d’Informatique Fondamentale, Université de Provence. Supported in part by the RFBR grant 06-01-00122-a

¹In the case of finite sets any such program is called an *implicit description of A* , as opposed to *explicit description of A* when M is required to halt after having printed the last element of A .

with one infinite input tape and one infinite output tape as described above. For every infinite 0-1-sequence ω let $M(\omega)$ denote the set enumerated by M when ω is written on its input tape. For every $A \subset \mathbb{N}$ consider the probability

$$m_M(A) = \Pr[M(\omega) = A].$$

A theorem of de Leeuw, Moore, Shannon and Shapiro [2] states that if $m_M(A) > 0$ then A is enumerable.

The class of distributions of such form has a maximal one up to a multiplicative constant. In other words, there is a machine M_1 (called *optimal*) such that for every machine M there is a constant c such that

$$c \cdot m_{M_1}(A) \geq m_M(A)$$

for all $A \subset \mathbb{N}$. Fix any such M_1 and call $m(A) \stackrel{def}{=} m_{M_1}(A)$ the *a priori* probability of enumerating A . The a priori distribution thus depends on the choice of the optimal machine but this dependence is also weak: for any other optimal machine M_2 both ratios $m_{M_1}(A)/m_{M_2}(A)$ and $m_{M_2}(A)/m_{M_1}(A)$ are bounded by a constant not depending on A . Let $H(A)$ denote the negative binary logarithm of the a-priori probability of A : $H(A) = \lceil -\log m(A) \rceil$.

Comparing M_0 , the machine defining $I(A)$, with M_1 , the machine defining $m(A)$, we see that

$$H(A) = \lceil -\log m_{M_1}(A) \rceil \leq I_{M_1}(A) \leq I_{M_0}(A) + O(1) = I(A) + O(1)$$

for all A . Solovay [5] has proved that conversely $I(A) \leq 3H(A) + O(\log H(A))$ for all A , which can be viewed as a sharpening of de Leeuw et al.'s result.

Theorem 1 (Solovay). *There is a constant c such that for every set $A \subset \mathbb{N}$ we have $I(A) \leq 3H(A) + 2 \log H(A) + c$.*

It is unknown whether we can replace the constant 3 in this inequality by a smaller constant. In this paper, we show that the constant 3 can be replaced by the constant 2 for *finite* sets A .

Theorem 2. *There is a constant c such that for every finite set A we have $I(A) \leq 2H(A) + 2 \log H(A) + c$.*

The proof of Theorem 2 is basically a simplification of that of Theorem 1. Thus we first sketch the latter one and then present the former one, explaining the main difference between two proofs. Our proof of Theorem 1 keeps the main ideas of Solovay's proof but differs from it in many technical details.

First we introduce some terminology and notation. Let Ω stand for the set of all infinite binary sequences. We write $x \leq \omega$ if x is a finite prefix of an infinite 0-1-sequence ω . Let Ω_x denote the set of all $\omega \in \Omega$ with $x \leq \omega$.

We will consider Cantor topology on Ω . Its base open sets are all sets of the form Ω_x . We call a subset of Ω *finitely based* if it is a finite union of sets of the form Ω_x .

There is a natural one-to-one correspondence between the family of all subsets of \mathbb{N} and Ω . Each subset A of \mathbb{N} corresponds to its characteristic sequence, whose n th bit is 1 iff $n \in A$. In what follows we will identify subsets of Ω with their characteristic sequences. In particular, we will write $A \in \Omega_x$ and $x \leq A$ to indicate that x is a prefix of the characteristic sequence of A . The notation $A \subset B$ will be used for the inclusion relation.

Fix an optimal machine M defining the a priori distribution m . Let $M^t(\omega)$ stand for the set enumerated by M in t steps on input ω . For each $A \subset \mathbb{N}$ and t let

$$S^t(A) = \{\omega \mid M^t(\omega) = A\}.$$

The set $S^t(A)$ is finitely based and a code of $S^t(A)$ (a finite list of respective x 's) can be computed given t and A . (For all infinite A and all t the set $S^t(A)$ is empty.) Let μ denote the uniform measure on Ω and $m^t(A) = \mu(S^t(A))$. According to our agreement to identify sets with their characteristic sequences, we denote by $m^t(\Omega_x)$ the total m^t -measure of all sets whose characteristic sequence begins with x .

Note that $m^t(A)$ can both increase and decrease as t increases (if A is finite). Indeed, assume that $M^{t-1}(\omega) = A$ and on step t of the run on input ω the machine M writes a new element b on the output tape. Let x be the length- t prefix of ω . Then $S^t(A)$ is decremented by Ω_x , while $S^t(A \cup \{b\})$ is incremented by Ω_x on step t .

Proof of Theorem 1 (a sketch). Let T be a subset of $\{0, 1\}^*$. A limit point of T is an infinite 0-1-sequence ω having the following property: Every its prefix is a prefix of some string in T . By T^n we denote the set of all strings of length n in T . We say that an algorithm constructs a set T of binary strings if for any given n it prints the list of strings in T^n and then halts.

An essential part of the proof is an algorithm that for every k constructs a set of strings T_k and computes a sequence of natural numbers

$$t_0 < t_1 < t_2 < \dots$$

having the following properties:

$$m^{t_n}(\Omega_x) \geq 2^{-k-1} \text{ for all } n \text{ and all } x \in T_k^n, \quad (1)$$

$$\text{if } m(A) \geq 2^{-k} \text{ then } A \text{ is a limit point of } T_k. \quad (2)$$

Lemma 1. *There is an algorithm that given k and an auxiliary binary string λ_k of length $k + 1$ constructs a set T_k and computes an increasing sequence t_0, t_1, \dots having the properties (1) and (2).*

For the goal of this paper, we do not need the proof of this lemma. However, for the sake of completeness we present its proof in the Appendix.

So, assume Lemma 1. We are going to construct an algorithm that on input λ_k enumerates certain subsets C_1, \dots, C_N of \mathbb{N} , where $N = O(2^{2k})$, having the following property: Every limit point A of T_k (the tree constructed by the algorithm of Lemma 1) is among C_1, \dots, C_N .

To this end we need a computable strategy to win an infinite two person game defined in Martin's paper [4]. Let N, K be natural numbers. A configuration in this game is an N -tuple of finitely based subsets of Ω : $\langle Z_1, \dots, Z_N \rangle$. The initial configuration is $\langle \Omega, \dots, \Omega \rangle$. Player I on his turn plays a finitely based set Y with $\mu(Y) \geq 1/K$. If Y is disjoint with Z_i for all $i = 1, \dots, N$, player I wins. If not, player II chooses a Z_i intersecting Y and replaces Z_i by Y , and the game continues. Player II wins if he can prevent I to win as described above for the entire, infinitely long game.

Martin [4] has proved that if $N = K(K + 1)/2$ then player II has a computable winning strategy (uniformly in K). We use Martin's result for $K = 2^{k+1}$. (We present its proof in the Appendix.)

Algorithm. We make steps $n = 1, 2, \dots$. At the end of step n we will have a configuration $\langle Z_1, \dots, Z_N \rangle$ in Martin's game and the sets C_1, \dots, C_N enumerated so far. They will satisfy the following conditions:

$$C_i \subset \{1, \dots, n\} \text{ for all } i \leq N; \quad (3)$$

$$C_i \subset M^{t_n}(\omega) \text{ for all } \omega \in Z_i \text{ and all } i \leq N; \quad (4)$$

$$\text{every } x \in T_k^n \text{ is a prefix of } C_i \text{ for some } i \leq N; \quad (5)$$

$$\text{the configuration } \langle Z_1, \dots, Z_N \rangle \text{ was obtained by applying a computable} \quad (6)$$

$$\text{winning strategy of player II against a sequence of moves of player I.}$$

At the start $n = 0$, $C_i = \emptyset$ and $Z_i = \Omega$, and conditions (3), (4), (5) and (6) are straightforward.

Step n . At the beginning of step n the conditions (3), (4), (5) and (6) are true for $n - 1$. The conditions (3) and (4) for $n - 1$ imply conditions (3) and (4) for n . Thus all the conditions except (5) are true at the beginning of step n . The condition (5) however may become false for any of $x \in T_k^n$, as T_k^n and T_k^{n-1} can be unrelated.

How to restore condition (5)? First, using the algorithm of Lemma 1, find the list of T_k^n . Then pick any x from T_k^n , call it x_1 . Play $Y_1 = S^{t_n}(\Omega_{x_1})$ for player I in Martin's game. By condition (1) this is a legal move. Assume that the winning strategy of player II plays Z_i .

We want to add some elements to C_i to ensure $x_1 \leq C_i$. To this end we need to show that C_i is a subset of the set $X_1 = \{j \mid j\text{th bit of } x_1 \text{ is } 1\}$. We are given that there is $\omega \in Z_i$ which is in Y_1 , that is,

$$x_1 \leq M^{t_n}(\omega).$$

By condition (4) we have

$$C_i \subset M^{t_n}(\omega).$$

Thus C_i is a subset of a set whose characteristic sequence begins with x_1 , therefore $C_i \cap \{1, \dots, n\} \subset X_1$. By condition (3) this implies that thus $C_i \subset X_1$.

Update C_i by letting $C_i = X_1$. Thus condition (5) is fulfilled for $x = x_1$.

Replace Z_i by Y_1 . This replacement restores condition (4) (that might become false after changing C_i). Indeed, $x_1 \leq M^{t_n}(\omega)$ for all $\omega \in Y_1$ and C_i has just become equal to X_1 .

Next we pick another element x_2 from T_k^n and repeat the procedure for x_2 in place of x_1 . Note that the set Y_2 is disjoint with Y_1 , as $M^{t_n}(Y_2)$ and $M^{t_n}(Y_1)$ are equal to disjoint sets Ω_{x_2} and Ω_{x_1} . Thus the new move Z_i of player II is different from Y_1 and the condition (5) remains true for $x = x_1$. Repeating the procedure $|T_k^n|$ times we fulfil condition (5) for all $x \in T_k^n$. **End of Algorithm.**

We need to prove that for every limit point A of T_k there is i with $C_i = A$. For every prefix x of (the characteristic sequence of) A , at the end of some step $n \geq |x|$ the string x is a prefix of C_i for some i . That i may depend on x . However, as the number of possible i 's is finite, there is i such that every prefix x of A is a prefix of C_i at the end of some step $n \geq |x|$.

This obviously implies that $A \subset C_i$. To prove the converse inclusion, pick any $j \in C_i$ and assume that j was included in C_i on the step m (thus $m \geq j$). The length- m prefix of A is a prefix of C_i at the end of some step $n \geq m$. Thus j th bit of the characteristic function of A is 1.

Consider the machine that on every input ω beginning with

$$p = 0^{\log k} 1(\text{binary notation of } k)(\lambda_k)(\text{binary notation of } i)$$

scans p and then, running the Algorithm, enumerates the set C_i (and no other sets among C_1, \dots, C_N). For this machine M' it holds

$$I_{M'}(C_i) \leq 2 \log k + 1 + k + 1 + \log O(2^{2k})$$

and by universality

$$I(C_i) \leq I_{M'}(C_i) + O(1) \leq 3k + 2 \log k + O(1)$$

for all i .

Let A be any enumerable subset of \mathbb{N} and $k = H(A)$. By condition (2) there is i such that the set C_i enumerated by the Algorithm coincides with A . Thus we obtain

$$I(A) \leq 3H(A) + 2 \log H(A) + O(1). \quad \square$$

How can we improve Solovay's bound

$$I(A) \leq 3H(A) + 2 \log H(A) + O(1)?$$

We could try to improve the upper bound of N in Martin's game. However, Ageev [1] showed that the condition $N = \Omega(K^2)$ is necessary for player II to win. Another option would be to reduce the length of the auxiliary string λ_k in Lemma 1. We do not know if this is possible.

For finite sets we can simplify the above construction as follows. If A is finite and $m(A) \geq 2^{-k}$ then $m^t(A) \geq 2^{-k-1}$ for all large enough t (Lemma 2 below). Thus we do not need the algorithm of Lemma 1. The algorithm that enumerates sets C_1, \dots, C_N enforces, on step t , that every finite A with $m^t(A) \geq 2^{-k-1}$ be among C_1, \dots, C_N . Thus we get rid of the string λ_k , and the enumeration complexity of C_1, \dots, C_N is reduced by about k bits.

Proof of Theorem 2. We construct an algorithm that given k enumerates $N = O(2^{2k})$ sets C_1, \dots, C_N so that every finite set A with $m(A) \geq 2^{-k}$ coincides with C_i for some $i \leq N$. Just as in the proof of Theorem 1 this implies

$$I(A) \leq 2H(A) + 2 \log H(A) + O(1).$$

Algorithm. We make steps $t = 1, 2, \dots$. At the end of each step t we will have a configuration $\langle Z_1, \dots, Z_N \rangle$ in Martin's game for $K = 2^{k+1}$ and the sets C_1, \dots, C_N enumerated so far. They will satisfy the following conditions:

$$C_i \subset M^t(\omega) \text{ for all } \omega \in Z_i \text{ and all } i = 1, \dots, N, \quad (7)$$

$$\text{every finite } A \text{ with } m^t(A) \geq 2^{-k-1} \text{ is among } C_1, \dots, C_N, \quad (8)$$

$$\text{the configuration } \langle Z_1, \dots, Z_N \rangle \text{ was obtained by applying a computable} \quad (9)$$

$$\text{winning strategy of player II against a sequence of moves of player I.}$$

Step t . At the beginning of step t conditions (7), (8) and (9) are true for $t - 1$. Obviously, condition (7) remains valid for t in place of $t - 1$, and we need to restore condition (8).

First we find all sets A_1, \dots, A_s with

$$m^t(A_i) \geq 2^{-k-1}.$$

Then we play $Y_1 = S^t(A_1)$ for the player I in Martin's game. Let Z_i be the move of the computable winning strategy of player II. As Z_i intersects Y_1 , there is $\omega \in Z_i$ with $A_1 = M^t(\omega)$. By condition (7) we have

$$C_i \subset M^t(\omega) = A_1.$$

We update C_i by letting $C_i = A_1$ and replace Z_i by Y_1 . These changes enforce condition (8) for $A = A_1$ and do not break condition (7).

Then we repeat the procedure for A_2, \dots, A_s in place of A_1 . Note that $S^t(A_1), \dots, S^t(A_s)$ are pairwise disjoint. Therefore, the condition (8) for A_1, \dots, A_{i-1} will not be broken when we perform the procedure for A_i . **End of Algorithm.**

The correctness of the Algorithm is based on the following

Lemma 2. $m(A) = \lim_{t \rightarrow \infty} m^t(A)$ for every finite $A \subset \mathbb{N}$.

The proof of this lemma is an easy exercise in measure theory and is given in the Appendix.

If A is finite and $m(A) \geq 2^{-k}$ then by Lemma 2 for almost all t we have $m^t(A) \geq 2^{-k-1}$. Therefore there is i such that on infinitely many steps we have $C_i = A$. Since C_i can only increase on each step, A coincides with C_i starting from some step. \square

Acknowledgments. The author is sincerely grateful to Sergei Salnikov for writing down a preliminary version of the proof and for the anonymous referees for helpful suggestions.

References

- [1] M. Ageev, Martin's game: a lower bound for the number of sets, *Theor. Comput. Sci.* 289(1): 871-876 (2002)
- [2] K. de Leeuw, E. F. Moore, C. E. Shannon, and N. Shapiro, Computability by probabilistic machines, In: C. E. Shannon and J. McCarthy (Eds.), *Automata Studies*, Princeton University Press, Princeton, New Jersey, 1956, 183-212.
- [3] M. Li and P.M.B. Vitányi, *An Introduction to Kolmogorov Complexity and its Applications*, Springer-Verlag, New York, 2nd Edition, 1997.
- [4] D.A. Martin, Borel indeterminacy, *Ann. Math.* 102 (1978) 363-371.
- [5] R.M.Solovay, On Random R.E. Sets, In: A.I. Arruda, N.C.A. da Costa, R. Chaqui (Eds.), *Non-Classical Logics, Model Theory and Computability*, North-Holland, Amsterdam, 1977, pp. 283-307.

Appendix.

Proof of Lemma 2. The set $S^t(A)$ is the difference of two sets: $S_1^t = \{\omega \mid M(\omega) \text{ prints in at most } t \text{ steps all the elements of } A\}$ and $S_2^t = \{\omega \mid M_1(\omega) \text{ prints in at most } t \text{ steps all the elements of } A \text{ and an element of } \mathbb{N} \setminus A\}$. Let S_1^∞ be the union of all S_1^t and S_2^∞ the union of all S_2^t . As the uniform measure is continuous we have

$$\mu(S_1^\infty) = \lim_{t \rightarrow \infty} \mu(S_1^t), \quad \mu(S_2^\infty) = \lim_{t \rightarrow \infty} \mu(S_2^t),$$

and

$$\begin{aligned} m(A) &= \mu(S_1^\infty \setminus S_2^\infty) \\ &= \mu(S_1^\infty) - \mu(S_2^\infty) \\ &= \lim_{t \rightarrow \infty} \mu(S_1^t) - \lim_{t \rightarrow \infty} \mu(S_2^t) \\ &= \lim_{t \rightarrow \infty} (\mu(S_1^t) - \mu(S_2^t)) \\ &= \lim_{t \rightarrow \infty} \mu(S_1^t \setminus S_2^t) = \lim_{t \rightarrow \infty} \mu(S^t(A)) = \lim_{t \rightarrow \infty} m^t(A). \quad \square \end{aligned}$$

Proof of Lemma 1. Let A_1, \dots, A_r be all sets with $m(A_i) \geq 2^{-k}$ and let $\lambda = \sum_{i=1}^r m(A_i)$. Let λ_k be the rational number consisting of $k+1$ first binary digits of the number $\lambda - 2^{-k-1}$.

The number t_n and the list T_k^n of all strings of length n in T_k are defined recursively. Let t_n be the first $t > t_{n-1}$ such that

$$m^t(\Omega_{x_1} \cup \dots \cup \Omega_{x_s}) \geq \lambda_k,$$

where x_1, \dots, x_s are all binary strings of length n with

$$m^t(\Omega_{x_i}) \geq \max\{2^{-k-1}, 2^{-k} - 1/n\}.$$

Let $T_k^n = \{x_1, \dots, x_s\}$.

We have to prove first that such t exists. Let x_1^n, \dots, x_r^n stand for prefixes of length n of characteristic functions of A_1, \dots, A_r (some of them may coincide). For any finitely based set $S \subset \Omega$ we have

$$m(S) = \lim_{t \rightarrow \infty} m^t(S).$$

(This can be proved just as Lemma 2.) This implies that for all large enough t we have

$$m^t(\Omega_{x_1^n} \cup \dots \cup \Omega_{x_r^n}) \geq \lambda_k$$

We can pick t so large that we additionally have

$$m^t(\Omega_{x_i^n}) \geq \max\{2^{-k-1}, 2^{-k} - 1/n\}.$$

for all $i = 1, \dots, r$. Indeed, $m^t(\Omega_{x_i^n})$ tends to $m(\Omega_{x_i^n})$, which is at least $m(\Omega_{A_i}) \geq 2^{-k}$. Any such t qualifies all the requirements.

It remains to prove that every A with $m(A) \geq 2^{-k}$ is a limit point of T_k . Let B_1, \dots, B_m be all different limit points of T_k . It suffices to show that

$$m(\{B_1, \dots, B_m\}) \geq \lambda_k \quad \text{and} \quad m(B_i) \geq 2^{-k}$$

for all $i \leq m$. (Indeed, if A was not among B_1, \dots, B_m , then the m -measure of the set of all B 's with $m(B) \geq 2^{-k}$ would be at least $\lambda_k + 2^{-k} > \lambda$.)

Let us prove first that $m(B_i) \geq 2^{-k}$ for every $i \leq m$. Fix i and let z_i^n denote length- n prefix of B_i . As m is a continuous measure, we have

$$m(B_i) = \lim_{n \rightarrow \infty} m(\Omega_{z_i^n}).$$

Thus it suffices to prove that $m(\Omega_{z_i^n}) \geq 2^{-k}$ for all n . Fix n . As B_i is a limit point of T_k , for every $j \geq n$ there exists $l \geq j$ such that z_i^j is a prefix of some $x \in T_k^l$. This implies that

$$m^{t_l}(\Omega_{z_i^n}) \geq m^{t_l}(\Omega_{z_i^j}) \geq m^{t_l}(\Omega_x) \geq 2^{-k} - 1/l.$$

Thus for every $j \geq n$ there is $l \geq j$ with

$$m^{t_l}(\Omega_{z_i^n}) \geq 2^{-k} - 1/l,$$

which implies that $m(\Omega_{z_i^n}) \geq 2^{-k}$.

The inequality $m(\{B_1, \dots, B_m\}) \geq \lambda_k$ is proved in a similar way. Let S_n stand for the set of all ω whose length- n prefix is among of z_1^n, \dots, z_m^n :

$$S_n = \Omega_{z_1^n} \cup \dots \cup \Omega_{z_m^n}.$$

The m -measure of S_n tends to $m(\{B_1, \dots, B_m\})$, thus it suffices to show that

$$m(S_n) \geq \lambda_k$$

for all n .

Fix any n . For all large enough l every string x of length l in T_k is an extension of some string among z_1^n, \dots, z_m^n . (Otherwise, compactness arguments show that T_k has a limit point outside S_n , hence different from B_1, \dots, B_m). Fix any such l and let x_1, \dots, x_s denote all strings of length l in T_k . By construction we have

$$m^{t_l}(\Omega_{x_1} \cup \dots \cup \Omega_{x_s}) \geq \lambda_k$$

and therefore

$$m^{t_l}(S_n) \geq m^{t_l}(\Omega_{x_1} \cup \dots \cup \Omega_{x_s}) \geq \lambda_k.$$

Since this inequality holds for all large enough l and $m(S_n) = \lim_{t \rightarrow \infty} m^t(S_n)$, we are done. \square

How to win Martin's game? Without loss of generality we may assume that the measure of every move Y of player I is exactly $1/K$ and thus at any moment of the game the measure of all sets Z_1, \dots, Z_N is exactly $1/K$. (At the start of the game we will reduce Z_1, \dots, Z_N ; if we win the game with reduced Z_1, \dots, Z_N then we certainly win with original Z_1, \dots, Z_N .)

We assign to every set Z_i a natural number in the range $1, \dots, K$, called the *rank*, so that (1) for every $r \leq K$ there are exactly r sets of rank r and (2) all sets of the same rank are pair wise disjoint (thus sets of rank K form a partition of Ω).

The condition (2) implies that each move Y of player II intersects some set of rank K . On our next move we choose a set Z_i of lowest rank r that intersects Y , replace it by Y and assign the rank $r - 1$ to Y . The condition (2) is thus satisfied. Note that now there are r (pair wise disjoint sets) of rank $r - 1$ and $r - 1$ sets of rank r . Thus swapping sets of rank r and $r - 1$ restores the condition (1).