# ASYMPTOTIC ESTIMATE OF THE LENGTH OF A DIAGNOSTIC
# WORD FOR A FINITE AUTOMATON

I. K. Rystsov                                                    UDC 519.713.4

By a finite automaton we understand a finite complete automaton with an output as strictly defined in [1] where basic concepts of the theory of automata can be found. Gill [2] has defined the concept of a diagnostic word for an automaton and derived the upper bound for its length in the form $n^n$, where n is the number of states of the automaton. This upper bound has been subsequently repeatedly reduced, the best of all upper bounds being of the order $2^{n/2}$ [3]. The lower bound $2^{n/4}$ has been first obtained in [4]. The purpose of this work is to bridge the gap between the lower and upper bounds.

Let L(n) denote the length of the longest of all shortest diagnostic words for automata having n states and at least one diagnostic word. In the following we will show that for any positive number $\varepsilon$ and beginning with large enough n the following properties are satisfied:

$$3^{\frac{n}{6}(1-\varepsilon)} < L(n) < 3^{\frac{n}{6}(1+\varepsilon)}. \tag{1}$$

It is easy to see that this gives an asymptotic estimate of the form $\log_3 L(n) \sim n/6$, or more accurately

$$\lim_{n \to \infty} \frac{6 \cdot \log_3 L(n)}{n} = 1.$$

First let us introduce the concept of a partition transversal which somewhat differs from the concept of a transversal as used in [3]. Let U denote a finite nonempty set of m elements and $\eta$ a partition on this set or, which is the same, an equivalence relation on the set U.

Definition 1. The subset $U_1 \subseteq U$ is called a transversal of the partition $\eta$ if each class of the partition contains not more than one element of the subset $U_1$.

By tr$(\eta, j)$ we denote the number of all transversals of partition $\eta$ of magnitude j. Let the rank of partition $\eta$ be k and let $m_i$, $1 \le i \le k$, be the size of the i-th class; then, from Definition 1 it is easy to see that for $j \le k$

$$\text{tr}(\eta, j) = \sum_{1 \le i_1 < i_2 < \cdots < i_j} m_{i_1} m_{i_2} \dots m_{i_j}, \tag{2}$$

where the sum is taken over all combinations of k subscripts j at a time. In particular, for j = k we obtain the number of "complete" transversals:

$$\text{tr}(\eta, k) = \prod_{i=1}^{k} m_i.$$

The following lemma gives an upper bound for the number of complete transversals.

LEMMA 1. The inequality

$$\text{tr}(\eta, k) \le 3^{\frac{m}{3}},$$

is true for any partition $\eta$ having k classes on a set of m elements.

Proof. As can be easily verified, the inequality $z \le 3^{z/3}$ holds for any natural number $z \ge 1$. Denoting, further, by $m_1, m_2, \dots, m_k$ the sizes of the classes of partition $\eta$, we have

$$\text{tr}(\eta, k) = \prod_{i=1}^{k} m_i \le \prod_{i=1}^{k} 3^{\frac{m_i}{3}} = 3^{\frac{m}{3}}.$$

---

This proves the lemma.

In fact, the number of complete transversals can be estimated more precisely (see [3]) using the integral function f(m) which is defined for all natural $m \geq 2$ as follows:

$$f(m) = \begin{cases} 3^k, & m = 3k, \\ 3^{k-1} \cdot 4, & m = 3k + 1, \\ 3^k \cdot 2, & m = 3k + 2. \end{cases} \tag{3}$$

Since this is immaterial in our case, we have given the simpler proof. It should be only pointed out that definition (3) makes it clear that for all natural $m \geq 2$

$$3^{\frac{m-3}{3}} \leqslant f(m). \tag{4}$$

Lemma 1 provides an upper bound for the number of transversals in partition with magnitudes between h and k, where $h \leq k$. In fact, for any combination of j subscripts we have $m_{i_1} \ldots m_{i_j} \leq m_1 m_2 \ldots m_k$. Then, from (2) we have $\operatorname{tr}(\eta, j) \leq \operatorname{tr}(\eta, k) \cdot C_k^j$, where $C_k^j$ is the binomial coefficient. Hence and from Lemma 1 we have

$$\sum_{j=h}^{k} \operatorname{tr}(\eta, j) \leqslant \sum_{j=h}^{k} \operatorname{tr}(\eta, k) \cdot C_k^j \leqslant 3^{\frac{m}{3}} \sum_{j=h}^{k} C_k^j.$$

Thus, assuming that the subscript i is equal to $k - j$, we have the inequality

$$\sum_{j=h}^{k} \operatorname{tr}(\eta, j) \leqslant 3^{\frac{m}{3}} \sum_{l=k-h}^{0} C_k^{k-l} = 3^{\frac{m}{3}} \sum_{l=0}^{k-h} C_k^l. \tag{5}$$

Later we shall need another type of automata: partial automata. A partial automaton is a finite automaton without an output and defined by the triple $B = (U, X, \gamma)$, where U and X are finite nonempty sets of states and inputs, and the function $\gamma : U \times X \to U$ is a function of transitions, generally speaking not everywhere defined. An input word is defined as a sequence of input signals. The transition function is extended in the usual way to the set of input words [1]. It is assumed that the effect of an empty word is to turn any state into itself. An input word is said to be admissible for the state $u \in U$ if a transition from the state u under the action of the word p is defined; in this case the state $\gamma(u, p)$ is denoted by up. An input word p is said to be admissible for the subset of states $U_1$ if a transition is defined for all states $u \in U_1$ under the action of the word p; in this case $U_1 p$ denotes the subset $\{up \mid u \in U_1\}$. Words admissible for the entire set of states are said to be admissible for the automaton B or simply admissible. (In [3] such words were called allowable.) Note that an empty word is admissible for any partial automaton. The number of states in the subset $U_1$ is denoted by $|U_1|$.

For any admissible word p we can define a partition on the set of states of an automaton following the expression

$$\eta(p) = \{(u_1, u_2) \mid \gamma(u_1, p) = \gamma(u_2, p)\}.$$

Let $U_1, U_2, \ldots, U_k$ be classes of the partition $\eta(p)$; from the above definition it is obvious that in such a case the effect of the word p is to turn each class $U_i$ into one state $\{U_i p\} = \{v_i\}$, $1 \leq i \leq k$, all states $v_i$ being distinct. Hence, the number of states in the subset Up is equal to the rank of the partition $\eta(p)$. Note that if $U_1$ is not a transversal of the partition $\eta(p)$, we have the strict inequality $|U_1 p| < |U_1|$.

Definition 2. An admissible word p is said to be irredundant for the automaton B if for any word q admissible for the set Up we have the equality $|Upq| = |Up|$.

From this definition it is clear that the number of states in subset Up is independent of the choice of a particular irredundant word p but depends on the automaton B. This number is called the degree of compressibility of B and is denoted by g(B). Note also that irredundant words exist in any finite partial automaton, so that T(B) will denote the length of the shortest irredundant word in automaton B. Let us now define the function T(m) as follows:

$$T(m) = \max_B \{T(B) \mid B \in \mathfrak{B}_m\},$$

where $\mathfrak{B}_m$ is the set of all partial automata with m states. This definition is correct since the set $\mathfrak{B}_m$ can be assumed to be finite if we stipulate that different input signals cause different partial transitions on the set of states.

The importance of the function T(m) becomes clear from the following theorem which was proved in [3].

THEOREM 1. For all natural numbers n ≥ 6 we have the property

$$T\left(\left[\frac{n}{2}\right]\right) \leqslant L(n) \leqslant T\left(\left[\frac{n}{2}\right]\right) \cdot n^2,$$

where [y] denotes the integral part of the natural number y.

However, the estimate of the function T(m) is also of special interest as it is related directly to the estimate of the length of the synchronizing word in a partial automaton. Namely, an admissible word p is said to be synchronizing for the automaton B if |Up| = 1. Obviously, of g(B) = 1, any irredundant word will be synchronizing for the automaton B and vice versa. It can be easily shown that the function T(m) is maximum in an automaton B in which g(B) = 1. Thus, the function T(m) is equal to the length of the largest of all shortest synchronizing words for automata of $\mathfrak{B}_m$ which have at least one synchronizing word.

Let us now turn directly to finding estimates for the function T(m). Gill ([3], Theorem 2) obtained the inequality f(m) ≤ T(m), where m ≥ 3 and f(m) is defined according to expression (3). From this and from inequality (4) we have that for any real number ε > 0 and for large enough m the following inequalities will be realized:

$$3^{\frac{m}{3}(1-\varepsilon)} < 3^{\frac{m-3}{3}} \leqslant T(m). \tag{6}$$

Taking into account the obvious inequality $\left[\frac{n}{2}\right] \geqslant \frac{n}{2} - 1 = \frac{n}{2}\left(1 - \frac{2}{n}\right)$, from the bottom inequality of Theorem 1 and from the property (6) we obtain the lower bound in (1) for the function L(n).

To obtain the upper boundary one has to prove certain auxiliary assertions. Let $l$(p) denote the length of input word p, and recall that m is the number of states of automaton B.

LEMMA 2. Let p be a certain admissible word for automaton B which is not irredundant and let k = |Up|; then, for any natural number h, g(B) ≤ h < k, there can be found an admissible word q such that |Uq| ≤ h and the length of q does not exceed

$$3^{\frac{m}{3}} \cdot \left(\sum_{i=0}^{k-h} C_k^i\right) + (k-h+1) \cdot l(p).$$

Proof. The word q will be constructed step by step. First we apply to the automaton the word p and then $q_1$, where $q_1$ is the shortest admissible word for the subset Up such that either |$Upq_1$| < k or the subset $Upq_1$ is not a transversal of the partition η(p). In any one of these cases we have $m_1$ = |$Upq_1p$| < k = $m_0$. Note that such a word $q_1$ does necessarily exist as p is not irredundant and the length of $q_1$ is not loger than tr (η(p), $m_0$) since otherwise it will not be the shortest word having this property. Also note that $q_1$ can be empty if the subset Up is no longer a transversal of the partition η(p).

If $m_1$ ≤ h the construction process ends and its result is the word q = $pq_1p$. Otherwise the step is repeated but now with the word $pq_1p$. Namely, let $q_2$ be the shortest admissible word for the subset $Upq_1p$ such that either |$Upq_1pq_2$| < $m_1$ or the subset $Upq_1pq_2$ is not a transversal of the partition η(p). In any of these cases we have $m_2$ = |$Upq_1pq_2p$| < $m_1$. As in the preceding case, such a word $q_2$ does necessarily exist and its length does not exceed tr (η(p), $m_1$). If $m_2$ ≤ h, the process stops and its result is declared to be the word q = $pq_1pq_2p$. Otherwise, the construction step is repeated, etc.

Since $m_0$ > $m_1$ > $m_2$ > . . . > $m_r$, the construction process ends after a finite number of steps r, where r ≤ k − h, and its result is a word q of the form $pq_1pq_2 . . . q_rp$. It is seen from the construction that |Uq| = $m_r$ ≤ h, and the length of the word q is estimated by

$$l(q) = \sum_{i=1}^{r} l(q_i) + (r+1) \cdot l(p) \leqslant \sum_{j=h}^{k} \text{tr}(\eta(p), j) + (k-h+1) \cdot l(p).$$

Hence and from the inequality (5) follows the truth of Lemma 2.

Let us denote by ]y[ the smallest integral number greater than or equal to the real number y and consider the following lemma.

LEMMA 3. Let B be a partial automaton with m states; then, for all natural numbers r, 1 ≤ r ≤ m, and d, 1 ≤ d ≤ ]m/r[, there can be found an admissible word q such that |Uq| ≤ max (g(B), m − d · r) and its length

does not exceed

$$3^{\frac{m}{3}} \cdot (2+r)^{d-1} \cdot \left( \sum_{i=0}^{r} C_m^i \right).$$

_Proof._ The lemma is proved by induction on the number d. Let us fix a certain number r from the interval [1, m] and let d = 1. If an empty word is irredundant for the automaton B, the lemma is obviously satisfied. Otherwise the lemma is applied to an empty word and to the number max (g(B), m − r). We conclude then that there exists an admissible word q such that |Uq| ≤ max (g(B), m − r) and its length does not exceed

$$3^{\frac{m}{3}} \cdot \sum_{i=0}^{m-(m-r)} C_m^i = 3^{\frac{m}{3}} \sum_{i=0}^{r} C_m^i.$$

The lemma is thus proved also for this case.

Assume that Lemma 3 has been proved for all numbers smaller than or equal to d and let us prove it for d + 1. By assumption there exists a word p whose length does not exceed

$$l(p) \leqslant 3^{\frac{m}{3}} \cdot (2+r)^{d-1} \cdot \left( \sum_{i=0}^{r} C_m^i \right) \tag{7}$$

and for which |Up| ≤ max (g(B), m − d·r). If the word p is already irredundant the lemma is obviously true also for d + 1.

Let us now assume that the word p is not irredundant. Let k denote the number |Up| and h, max (g(B), m − d·r − r). By choice of p we have the inequalities g(B) < k ≤ m − d·r, so that the inequality k − h ≤ m − d·r − (m − dr − r) = r is satisfied. Hence, applying Lemma 2 to the word p and number h we conclude that there can be found an admissible word q such that |Uq| ≤ h and its length satisfies the inequality

$$l(q) \leqslant 3^{\frac{m}{3}} \left( \sum_{i=0}^{r} C_k^i \right) + (1+r) \cdot l(p). \tag{8}$$

Then from the inequalities (7) and (8) and the condition k ≤ m follows

$$l(q) \leqslant 3^{\frac{m}{3}} \left( \sum_{i=0}^{r} C_m^i \right) + (1+r) \cdot 3^{\frac{m}{3}} (2+r)^{d-1} \left( \sum_{i=0}^{r} C_m^i \right).$$

Hence, considering the obvious inequality $1 \leq (2+r)^{d-1}$, we have

$$l(q) \leqslant 3^{\frac{m}{3}} \cdot (2+r)^{d-1} \cdot \left( \sum_{i=0}^{r} C_m^i \right) + (1+r) \cdot 3^{\frac{m}{3}} \cdot (2+r)^{d-1} \left( \sum_{i=0}^{r} C_m^i \right).$$

Factoring out the common term, we obtain the required constraint on the length of the word q. This proves the lemma.

COROLLARY 1. In any partial automaton B with m states, for any natural number r, 1 ≤ r ≤ m there is a dead-end numberwh ose length is not greater than the number

$$3^{\frac{m}{3}} \cdot (2+r)^{\frac{m}{r}} \cdot \left( \sum_{i=0}^{r} C_m^i \right).$$

_Proof._ Make in Lemma 3 d equal to ]m/r[. Then, by virtue of the inequalities m − ]m/r[·r ≤ 0 < g(B), we conclude that the word q, whose existence is asserted in Lemma 3, is irredundant for the automaton B. To prove the corollary it is now only necessary to note that ]m/r[ − 1 ≤ m/r.

We can now turn to finding an upper bound for T(m).

THEOREM 2. For any real number ε > 0 there can be found a natural number $m_0$ such that for all natural numbers m ≥ $m_0$ we have the following inequality:

$$T(m) \leqslant 3^{\frac{m}{3}(1+\varepsilon)}.$$

_Proof._ Let there be given a certain positive real number ε. For a sufficiently large natural number r we have

$$(2+r)^{\frac{1}{r}} \leqslant 3^{\frac{\varepsilon}{6}}. \tag{9}$$

In fact, it is only necessary to select r so that the inequality

$$\frac{\log(2+r)}{r} \leqslant \frac{\varepsilon}{6},$$

is satisfied; this is always possible since the left side of the inequality approaches zero when r tends to infinity. Let us fix a certain number $r_0$ so that the inequality (9) is satisfied. With $r_0$ fixed, the expression $\sum_{i=0}^{r_0} C_m^i$ becomes a polynomial of degree $r_0$ of the variable m. However, since a polynomial increases slower than an exponential function, for a sufficiently large m, say $m \geq m_1$, where $m_1$ depends on $\varepsilon$ and $r_0$, we have the inequality

$$\sum_{i=0}^{r_0} C_m^i \leqslant 3^{\frac{m \cdot \varepsilon}{6}}. \tag{10}$$

Let us assume now that m is an arbitrary natural number greater than $m_0 = \max(r_0, m_1)$. Let us take an automaton B with m states such that $T(m) = T(B)$. Then, according to Corollary 1, we have

$$T(m) = T(b) \leqslant 3^{\frac{m}{3}}(2+r_0)^{\frac{m}{r_0}}\left(\sum_{i=0}^{r_0} C_m^i\right).$$

Hence and from the inequalities (9) and (10) follows the assertion of Theorem 2.

This proves the theorem.

As noted before, a power function increases more slowly than an exponential function; this means that for any given $\varepsilon > 0$ and a sufficiently large n we have the following inequality

$$n^2 \leqslant 3^{\frac{n \cdot \varepsilon}{6}}.$$

Hence and from Theorems 1 and 2 we obtain the upper bound for the function L(n) in property (1).

From Theorem 2 and inequality (6) we get that for any given $\varepsilon > 0$ and sufficiently large m the following inequalities are satisfied:

$$3^{\frac{m}{3}(1-\varepsilon)} \leqslant T(m) \leqslant 3^{\frac{m}{3}(1+\varepsilon)}. \tag{11}$$

These inequalities indicate that $\log_3 T(m)$ is asymptotically equal to m/3. In fact, taking logarithms for the base 3 of both sides of inequality (11) we obtain

$$\frac{m}{3}(1-\varepsilon) \leqslant \log_3 T(m) \leqslant \frac{m}{3}(1+\varepsilon).$$

Dividing both sides of these inequalities by m/3 we have

$$1-\varepsilon \leqslant \frac{3 \cdot \log_3 T(m)}{m} \leqslant 1+\varepsilon.$$

Hence, in view of the fact that $\varepsilon$ was arbitrarily selected, we obtain an asymptotic estimate for the logarithm of the function T(m): $\lim_{m \to \infty} \frac{3 \cdot \log_3 T(m)}{m} = 1$. An asymptotic estimate of the logarithm of the function L(n) can be obtained similarly.

## LITERATURE CITED

1. V. M. Glushkov, "Abstract theory of automata," Usp. Mat. Nauk, 16, No. 5, 3 (1961).
2. A. Gill, Introduction to the Theory of Finite Automata [Russian translation], Nauka, Moscow (1966).
3. I. K. Rystsov, "Estimate of the length of a diagnostic word for finite automata," Kibernetika, No. 6, 40 (1978).
4. M. N. Sokolovskii, "Estimate of the length of a diagnostic word for finite automata," Kibernetika, No. 2, 16 (1976).