

## Exponential Bounds for the Running Time of a Selection Algorithm

LUC DEVROYE

*School of Computer Science, McGill University, 805 Sherbrooke Street West,  
Montreal, Quebec H3A 2K6, Canada*

Received March 20, 1982; revised May 4, 1983

Hoare's selection algorithm for finding the  $k$ th-largest element in a set of  $n$  elements is shown to use  $C$  comparisons where

- (i)  $E(C^p) \leq A_p n^p$  for some constant  $A_p > 0$  and all  $p \geq 1$ ;
- (ii)  $P(C/n \geq u) \leq (\frac{3}{4})^{u(1+o(1))}$  as  $u \rightarrow \infty$ .

Exact values for the " $A_p$ " and " $o(1)$ " terms are given.

### 1. INTRODUCTION

Hoare [7], Aho *et al.* [1, pp. 101-102] and Horowitz and Sahni [9] all consider the following algorithm (with minor modifications) for finding the  $k$ th-smallest element in a set  $S$  of  $n$  elements ( $1 \leq k \leq n$ ):

```

procedure FIND ( $k, S$ )
  if  $|S| = 1$  then return the single element in  $S$ 
    else begin choose an element  $a$  randomly from  $S$ ;
      let  $S_1, S_2$  and  $S_3$  be the sequences of elements in  $S$  less than,
      equal to, and greater than  $a$ , respectively;
      if  $|S_1| \geq k$ , then return FIND ( $k, S_1$ )
      else if  $|S_1| + |S_2| \geq k$  then return  $a$ 
      else return FIND ( $k - |S_1| - |S_2|, S_3$ )
    end.
    
```

A nonrecursive version of this algorithm is of course easy to find. The work done here can be measured by the number of comparisons between elements (these occur only in the step in which  $S$  is split into  $S_1, S_2$  and  $S_3$ ). It is known that this algorithm requires  $\Omega(n^2)$  comparisons in the worst case. The algorithm runs in average time  $O(n)$  (Aho *et al.* [1]). In fact, Knuth [10] has shown that the average number of comparisons is at most

$$2((n+1)H_n - (n+3-i)H_{n-i+1} - (i+2)H_i + n + 3)$$

1

0022-0000/84 \$3.00

Copyright © 1984 by Academic Press, Inc.  
All rights of reproduction in any form reserved.

where  $H_n = \sum_{1 \leq j \leq n} (1/j)$ . Thus, for  $k = \sqrt{n/2}$ , we obtain the bound  $2(1 + \ln(2))n + o(n) \leq 3.39n + o(n)$ . For  $O(n)$  worst-case selection algorithms, see Blum *et al.* [2] or Schonhage *et al.* [11].

In this paper we give probabilistic bounds for the upper tail of  $C$ , the number of comparisons used by FIND. The bounds are reasonably tight, but more importantly, the exponential nature of the bounds shows that deviations from linearity are extremely unlikely. We do not want to challenge the fact that the algorithm of Floyd and Rivest [5, 6] is faster on the average than FIND (it was shown there that the expected number of comparisons is  $n + \min(k, n - k) + O(n^{1/2})$ ). The analysis given here for FIND is ad hoc, and therefore not directly extendible to the Floyd–Rivest algorithm.

*Result 1.* There exists a random variable  $T$ , independent of  $n$  and  $k$ , such that

$$C < nT$$

where  $<$  denotes “is stochastically smaller than,” i.e.,

$$P(C \geq u) \leq P(nT \geq u), \quad \text{all } u.$$

The random variable  $T$  satisfies  $E(T^p) < \infty$  for all  $p \geq 1$ .

*Result 2.*

$$\begin{aligned} E(C) &\leq 4n; \\ E(C^p) &\leq A_p n^p, \quad \text{all integer } p \geq 1, \end{aligned}$$

where

$$A_p = \frac{16}{3} \frac{p!}{\ln^{p-1}(\frac{4}{3})}.$$

*Result 3.*

$$P(C/n \geq u) \leq (1 + Au) e^{(\frac{3}{4})^u}, \quad u \geq 1/\ln(\frac{4}{3}),$$

and

$$P(C/n \geq u) \leq (\frac{3}{4})^{(\sqrt{u} - \sqrt{16/3})^2}, \quad u \geq 16/3,$$

where

$$A = (16/3) \ln^2(\frac{4}{3}).$$

## 2. ANALYSIS

We can and will assume that all elements in  $S$  are distinct. We claim that  $C$  is stochastically smaller than the outcome of the following algorithm.  $S$ ,  $k$  and  $n$  are as defined in the Introduction.

$l \leftarrow 0, r \leftarrow n + 1, C \leftarrow 0.$  ( $C$  will be the outcome of the algorithm.)  
 while  $r > l + 1$  do  
     begin generate  $N$  uniformly and at random in  $\{l + 1, \dots, r - 1\}$ ;  
          $C \leftarrow C + (r - l - 2)$ ;  
         if  $N < k$ , then  $(l, r) \leftarrow (N, r)$  else  $(l, r) \leftarrow (l, N)$   
 end.

Thus we can define  $C$  as the outcome of this algorithm, since we are only interested in upper bounds for  $C$ . In our proof, we will construct a probability space in the following manner. Let  $(U_1, V_1), (U_2, V_2), \dots$  be a sequence of independent uniform  $[0, 1]^2$  random vectors. We will use the notation  $(l_i, r_i)$  for the values of  $(l, r)$  in the  $i$ th iteration. In particular,  $(l_0, r_0) = (0, n + 1)$ . Our construction is such that the distribution of  $(l_i, r_i)$  is completely determined by  $(U_j, V_j), j \leq i$ .

Let  $(l_{i-1}, r_{i-1})$  be given. Then  $(l_i, r_i)$  is determined as follows:

$$(l_i, r_i) = \begin{cases} (l_{i-1}, r_{i-1} - 1 - \underline{(r_{i-1} - k)U_i}) & \text{if } V_i < p_{i-1} = (r_{i-1} - k)/(r_{i-1} - l_{i-1} - 1) \\ & \text{(an event that we shall call } A_i\text{);} \\ (l_{i-1} + 1 + \underline{(k - 1 - l_{i-1})U_i}, r_{i-1}) & \text{otherwise.} \end{cases}$$

Thus, on  $A_i$ ,  $r_i$  is uniformly distributed on  $\{k, \dots, r_{i-1} - 1\}$ , and on its complement,  $l_i$  is uniformly distributed on  $\{l_{i-1} + 1, \dots, k - 1\}$ . Thus, on  $A_i$ ,

$$\begin{aligned} r_i - l_i - 2 &= r_{i-1} - l_{i-1} - 3 - \underline{(r_{i-1} - k)U_i} \\ &\leq r_{i-1} - l_{i-1} - 2 - (r_{i-1} - k)U_i \\ &= r_{i-1} - l_{i-1} - 2 - p_{i-1}U_i(r_{i-1} - l_{i-1} - 1), \end{aligned}$$

and on  $A_i^c$ , the complement of  $A_i$ ,

$$\begin{aligned} r_i - l_i - 2 &= r_{i-1} - l_{i-1} - 3 - \underline{(k - 1 - l_{i-1})U_i} \\ &\leq r_{i-1} - l_{i-1} - 2 - (k - 1 - l_{i-1})U_i \\ &= r_{i-1} - l_{i-1} - 2 - (1 - p_{i-1})U_i(r_{i-1} - l_{i-1} - 1). \end{aligned}$$

Combining this, and using  $I$  to denote the indicator function of an event, gives

$$\begin{aligned} r_i - l_i - 2 &\leq (r_{i-1} - l_{i-1} - 2)(1 - U_i(p_{i-1}I_{A_i} + (1 - p_{i-1})I_{A_i^c})) \\ &= (r_{i-1} - l_{i-1} - 2)W_i \quad \text{(definition of } W_i\text{).} \end{aligned} \tag{1}$$

Inequality (1) is the starting point for all further analysis. Clearly, by recursion,

$$r_i - l_i - 2 \leq (n - 1) \prod_{j=1}^i W_j \leq n \prod_{j=1}^i W_j \tag{2}$$

and

$$C \leq n \left( 1 + \sum_{i=1}^{\infty} \prod_{j=1}^i W_j \right). \quad (3)$$

We will use two Lemmas.

LEMMA 1. *If  $U$  is a uniform  $[0, 1]$  random variable, then  $E((1 - U/2)^p) < 2/(p + 1)$ , all  $p \geq 1$ .*

$$\text{Proof. } E((1 - U/2)^p) = \int_0^1 (1 - u/2)^p du = (2/(p + 1))(1 - 2^{-(p+1)}) < 2/(p + 1).$$

LEMMA 2. *Let  $W_1, W_2, \dots$  be a sequence of independent identically distributed nonnegative random variables with  $p$ th moment  $E(W_1^p) = \mu < 1$ , and let  $X = 1 + \sum_{j=1}^{\infty} \prod_{i=1}^j W_i$ . For  $p \geq 1$ ,*

$$E(X^p) \leq 1/(1 - \mu^{1/p})^p. \quad (4)$$

For  $p = 1$ , equality is achieved in (4).

*Proof.* Whenever we have a random variable  $X$  that can be written as  $\sum_{i=0}^{\infty} X_i$ , then for all  $\lambda \in (0, 1)$ ,

$$X = \sum_{i=0}^{\infty} \lambda^i (1 - \lambda) \frac{X_i}{\lambda^i (1 - \lambda)}$$

so that by Jensen's inequality,

$$X^p \leq \sum_{i=0}^{\infty} \lambda^i (1 - \lambda) \left( \frac{X_i}{\lambda^i (1 - \lambda)} \right)^p = (1 - \lambda)^{1-p} \sum_{i=0}^{\infty} \lambda^{i(1-p)} X_i^p.$$

If we replace  $X_i$  by 1 for  $i = 0$  and by  $W_1 W_2 \dots W_i$  for  $i \neq 0$ , and if we note that  $E(X_i^p) = \mu^i$ , then

$$E(X^p) \leq \sum_{i=0}^{\infty} (\mu/\lambda^{p-1})^i (1 - \lambda)^{1-p} = (1 - \lambda)^{1-p} / (1 - \mu/\lambda^{p-1}), \quad \mu < \lambda^{p-1}.$$

The last expression is minimal for  $\lambda = \mu^{1/p}$ . Resubstitution gives the bound

$$(1 - \mu^{1/p})^{1-p} / (1 - \mu^{1-(p-1)/p}) = (1 - \mu^{1/p})^{-p}.$$

*Proof of Result 1.* It is clear that  $p_{i-1} I_{A_i} + (1 - p_{i-1}) I_{A_i^c} > \frac{1}{2} Z_i$ , where  $Z_i$  is Bernoulli with parameter  $\frac{1}{2}$  (note that  $A_i = [V_i < p_{i-1}]$ ). Thus, if  $Z_1, Z_2, \dots$  are independent Bernoulli ( $\frac{1}{2}$ ) random variables,

$$C < n \left( 1 + \sum_{i=1}^{\infty} \prod_{j=1}^i W_j^* \right) = nT \quad (\text{definition of } T) \quad (5)$$

where  $W_j^* = 1 - \frac{1}{2}Z_j U_j$ . Thus, all  $W_j^*$ 's are independent and identically distributed. Also,  $E(T) = 1 + \sum_{i=1}^{\infty} \prod_{j=1}^i E(W_j^*) = \sum_{i=0}^{\infty} (\frac{7}{8})^i = 8 < \infty$  so that the right-hand side of (5) is indeed almost surely finite. By Lemma 1,

$$E(W_1^{*p}) = \frac{1}{2} \left( 1 + E \left( 1 - \frac{U_1}{2} \right)^p \right) \leq \frac{1}{2} \left( 1 + \frac{2}{p+1} \right) = \frac{1}{2} \frac{p+3}{p+1}.$$

Thus, by Lemma 2, for  $p > 1$ ,

$$E(T^p) \leq 1 / \left( 1 - \left( \frac{1}{2} \frac{p+3}{p+1} \right)^{1/p} \right)^p < \infty.$$

*Remark 1.* The stochastic majorization used in this proof is sloppy. It gives the crude estimate  $E(C) \leq 8n$ . With the more refined majorization

$$p_{i-1} I_{A_i} + (1 - p_{i-1}) I_{A_i^c} > \frac{1}{2} (Z_i + (1 - Z_i) V_i)$$

where  $V_i$  is uniform  $[0, 1]$  and independent of  $Z_i$ , and with  $W_j^* = 1 - \frac{1}{2}(Z_j + (1 - Z_j)V_j)U_j$  in (5), we obtain the slightly sharper result

$$\frac{E(C)}{n} \leq \sum_{i=0}^{\infty} \left( \frac{13}{16} \right)^i = \frac{16}{3}.$$

*Proof of Result 2.* We take (3) as our starting point, and let  $\mathfrak{F}_i$  be the  $\sigma$ -algebra generated by  $(U_1, V_1), \dots, (U_i, V_i)$ .  $\mathfrak{F}_0$  is the  $\sigma$ -algebra consisting of the empty set and its complement. By well-known properties of conditional expectations (see, e.g., Chow and Teicher [4]),

$$E(W_i | \mathfrak{F}_{i-1}) = 1 - \frac{1}{2}(p_{i-1}^2 + (1 - p_{i-1})^2) \leq 1 - \frac{1}{4} = \frac{3}{4},$$

and

$$E \left( \prod_{i=1}^j W_i \right) = E \left( \prod_{i=1}^j E(W_i | \mathfrak{F}_{i-1}) \right) \leq \left( \frac{3}{4} \right)^j, \quad j \geq 1.$$

Since obviously  $0 \leq W_i \leq 1$ , we have for  $p \geq 1$ ,  $p$  integer,

$$\begin{aligned} \frac{E(C^p)}{n^p} &\leq E \left( \left( 1 + \sum_{i=1}^{\infty} \prod_{j=1}^i W_j \right)^p \right) = \sum_{j_1, \dots, j_p \in \{0, 1, 2, \dots\}^p} E \left( \prod_{m=1}^p \prod_{i < j_m} W_i \right) \\ &\leq \sum_{j_1, \dots, j_p \in \{0, 1, 2, \dots\}^p} E \left( \prod_{i < \max(j_1, \dots, j_p)} W_i \right) \\ &\leq \sum_{j_1, \dots, j_p \in \{0, 1, 2, \dots\}^p} \left( \frac{3}{4} \right)^{\max(j_1, \dots, j_p)} \\ &\leq p \sum_{j=0}^{\infty} (j+1)^{p-1} \left( \frac{3}{4} \right)^j = \frac{16}{3} p \sum_{j=0}^{\infty} j^{p-1} \frac{1}{4} \left( \frac{3}{4} \right)^j = \frac{16}{3} p E(X^{p-1}) \quad (6) \end{aligned}$$

where  $X$  is geometrically distributed:  $P(X = j) = \frac{1}{4}(\frac{3}{4})^j$ ,  $j \geq 0$ . We note that  $X$  is distributed as the integer part of  $X^*/\ln(\frac{4}{3})$ , where  $X^*$  is exponentially distributed (i.e., has density  $e^{-x}$  on  $[0, \infty)$ ). Thus (6) is bounded from above by

$$\frac{16}{3} p E(X^{*p-1})/\ln^{p-1} \left( \frac{4}{3} \right) = \frac{16}{3} p(p-1)!/\ln^{p-1} \left( \frac{4}{3} \right) = \frac{16}{3} p!/\ln^{p-1} \left( \frac{4}{3} \right).$$

The well-known result  $E(C) \leq 4n$  follows easily:

$$\frac{E(C)}{n} \leq 1 + \sum_{j=1}^{\infty} \prod_{i=1}^j E(W_i) \leq \sum_{j=0}^{\infty} \left( \frac{3}{4} \right)^j = 4.$$

*Proof of Result 3.* We start from Result 2. Let  $t$  be a real number in  $(0, \ln(\frac{4}{3}))$ , and let  $T$  be  $C/n$ . By Result 2,

$$E(e^{tT}) = \sum_{i=0}^{\infty} \frac{t^i}{i!} E(T^i) \leq 1 + \frac{16}{3} \sum_{i=1}^{\infty} \frac{t^i}{\ln^{i-1}(\frac{4}{3})} = 1 + \frac{16t}{3} \left( 1 - \frac{t}{\ln(\frac{4}{3})} \right)^{-1}. \quad (7)$$

Thus, by the Bernstein–Chernoff bounding method (see Chernoff [3] or Hoeffding [8]),

$$P(T \geq u) \leq E(e^{tT}) e^{-tu} \\ \leq \left( 1 + \frac{16t}{3} \left( 1 - t/\ln \left( \frac{4}{3} \right) \right)^{-1} \right) e^{-tu}, \quad 0 < t < \ln \left( \frac{4}{3} \right). \quad (8)$$

Result 3 now follows by choosing  $t$  carefully. For the first inequality, we take a positive number  $c$ , and assume that  $u > c/\ln(\frac{4}{3})$ ,  $t = \ln(\frac{4}{3}) - c/u$ . The last expression is not greater than

$$(1 + au/c) e^{c(\frac{3}{4})^u}$$

where  $a = 16 \ln^2(\frac{4}{3})/3$ . Considered as a function of  $c$ , the latter expression is minimal when  $c^2 + auc - au = 0$ , i.e., when  $c = (au/2)(\sqrt{1 + 4/au} - 1) \sim 1$  as  $au \rightarrow \infty$ . Thus the value  $c = 1$  is best for large  $u$ . This leads to the upper bound

$$(1 + au) e(\frac{3}{4})^u, \quad \text{valid for } u > 1/\ln(\frac{4}{3}).$$

For the second inequality of result 3, we apply the inequality  $1 + u \leq e^u$  to (8), and obtain the inequality

$$P(T \geq u) \leq \exp(-tu + (16t/3)(1 - t/\ln(\frac{4}{3}))^{-1}), \quad 0 < t < \ln(\frac{4}{3}), \quad (9)$$

which has the form  $\exp(-tu + at/(1 - bt))$ . Such an expression is minimal when  $t = (1 - \sqrt{a/u})(1/b)$ . Replacement of this value of  $t$  in (9) shows that

$$P(T \geq u) \leq \exp(-(\sqrt{u} - \sqrt{a})^2/b)$$

where  $a = 16/3$  and  $b = 1/\ln(\frac{4}{3})$ . The last inequality is valid for all  $u \geq a$ . This concludes the proof of Result 3.

## REFERENCES

1. A. V. AHO, J. E. HOPCROFT, AND J. D. ULLMAN, "The Design and Analysis of Computer Algorithms," Addison-Wesley, Reading, Mass., 1974.
2. M. BLUM, R. W. FLOYD, V. PRATT, R. L. RIVEST, AND R. E. TARJAN, Time bounds for selection, *J. Comput. System Sci.* **7** (1973), 448-461.
3. H. CHERNOFF, A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations, *Ann. Math. Statist.* **23** (1952), 493-507.
4. Y. S. CHOW AND H. TEICHER, "Probability Theory," Springer-Verlag, New York/Berlin, 1978.
5. R. W. FLOYD AND R. L. RIVEST, Expected time bounds for selection, *Comm. ACM* **18** (1975), 165-172.
6. R. W. FLOYD AND R. L. RIVEST, Algorithm 489, *Comm. ACM* **18** (1975), 173.
7. C. A. R. HOARE, Find (algorithm 65), *Comm. ACM* **4** (1961), 321-322.
8. W. HOEFFDING, Probability inequalities for sums of bounded random variables, *J. Amer. Statist. Assoc.* **58** (1963), 13-30.
9. E. HOROWITZ AND S. SAHNI, "Fundamentals of Computer Algorithms," Computer Science Press, Potomac, Md, 1978.
10. D. E. KNUTH, "Mathematical Analysis of Algorithms," Computer Science Dept. Report STAN-CS-71-206, Stanford University, 1971.
11. A. SCHONHAGE, M. PATERSON, AND N. PIPPENGER, Finding the median, *J. Comput. System Sci.* **13** (1976), 184-199.