

Increasing the Efficiency of Quicksort

M. H. VAN EMDEN

Mathematical Centre, Amsterdam, The Netherlands

A method is presented for the analysis of various generalizations of quicksort. The average asymptotic number of comparisons needed is shown to be $\alpha n \log_2(n)$. A formula is derived expressing α in terms of the probability distribution of the "bound" of a partition. This formula assumes a particularly simple form for a generalization already considered by Hoare, namely, choice of the bound as median of a random sample. The main contribution of this paper is another generalization of quicksort, which uses a *bounding interval* instead of a single element as bound. This generalization turns out to be easy to implement in a computer program. A numerical approximation shows that $\alpha = 1.140$ for this version of quicksort compared with 1.386 for the original. This implies a decrease in number of comparisons of 18 percent; actual tests showed about 15 percent saving in computing time.

KEYWORDS AND PHRASES: sorting, quicksort, information content, entropy, distribution of median

CR CATEGORIES: 3.73, 4.49, 5.31, 5.6

1. Quicksort

By sorting a sequence we mean arranging a sequence of numbers into nondecreasing order. We shall frequently refer to the *rank* of an element, which means the index of the place it occupies in the sorted sequence; the rank of a particular element, say x , will be denoted by $r(x)$. The sequence to be sorted will be referred to as $a[1], \dots, a[n]$ or as $a[1:n]$. Quicksort (which is due to C. A. R. Hoare—see [4, 5]) is the obvious choice for sorting a sequence that can be contained within a random access memory because it combines efficiency with the advantage that, apart from the sequence itself, the necessary additional storage is proportional to $\log_2(n)$.

The principle of quicksort may be described as follows. Take any real number y (let us refer to this as the *bound*).

This paper is related to an algorithm of the same title, which will be published in a later issue of *Communications of the ACM* [9].

Suppose that there are integers p and q such that

$$\left. \begin{array}{l} 0 \leq p < q \leq n + 1, \\ \text{among } a[1:p] \text{ none is greater than } y, \\ \text{among } a[q:n] \text{ none is smaller than } y. \end{array} \right\} \quad (1)$$

Suppose now that p and q are any pair of integers satisfying all three of these conditions and also such that $(q - p)$ is minimal. If $p + 1 = q$, then we are ready; otherwise we have $a[p + 1] > y$ and $a[q - 1] < y$. These elements are interchanged, and again p is increased and q is decreased as much as possible. This results in a change of at least 1 in p or q so that, continuing in this way, we at last find that $p + 1 = q$. Let the final value of p be r . A *partition* is now completed, which means that of $a[1:r]$ none is greater, and of $a[r + 1:n]$ none is smaller than y .

The problem of sorting $a[1:n]$ is now reduced to sorting $a[1:r]$ and $a[r + 1:n]$ separately. Thus the length of the sequence to be sorted is successively reduced until only sequences of length 1 or 2 are left. We shall see that the efficiency of quicksort depends on the way of selecting y . For the theoretical analysis of the number of comparisons, it will not be necessary to specify how y is selected; only the existence of a probability distribution will be postulated for the resulting r .

2. Average Number of Comparisons

Hoare [5] showed that the average number of comparisons done by quicksort is, asymptotically for large n , equal to $2n \ln(n)$. He also showed, by means of an information-theoretic argument, that the minimum number of comparisons is $n \log_2(n)$. This minimum would be achieved if the bound is always in the middle.

This section derives, by means of a similar argument, that the average number of comparisons is asymptotically equal to $\alpha n \log_2(n)$, where α is independent of n and we express α in terms of the probability distribution of the bound. Applications of this formula to the distributions considered by Hoare yield his results.

The number of comparisons required by a particular sorting algorithm may vary greatly for sequences of the same length; so, if we are to compare the performance of different sorting algorithms, we must define some sort of average. Suppose that no two elements are equal to each other; in that case not the elements themselves are important, but only their rank; hence a sequence may be regarded as a permutation of the integers $1, \dots, n$. In the sequel we mean by "average" the average over the ensemble of all permutations of $1, \dots, n$ where each occurs with the same probability.

An information-theoretic argument yields a formula for the average number of comparisons necessary to sort a sequence of length n . We can view the process of sorting as

one of collecting information about the particular permutation initially presented, because we would be able to reconstruct this permutation working backward from the sorted sequence if a record were kept of interchanges effected. Therefore, we must gain an amount of information equal to the uncertainty inherent in the random drawing of a permutation.

We suppose that, before the partition, any of the permutations of $1, \dots, n$ are equally probable. According to Shannon's theory of information [7], the uncertainty in this situation equals the entropy of the discrete probability distribution $\{p_1, \dots, p_n\}$ where $p_i = 1/n!$ for $i = 1, \dots, n!$:

$$-\sum_{i=1}^{n!} p_i \log_2 (p_i) = \log_2 (n!) \text{ bits.}$$

If, before the partition, $a[1:n]$ contains any of the possible permutations of $1, \dots, n$ with probability $1/n!$, then afterward the equivalent statement holds for $a[1:r]$ and $a[r+1:n]$. This may be verified as follows. The algorithm described in Section 1 may also be described as the successive random drawing without replacement from an urn initially containing balls numbered $1, \dots, n$. Suppose that $p + n - q + 1$ balls have been drawn already. Another ball is drawn, if available, and it is placed in $a[p+1]$ if smaller than y and in $a[q-1]$ otherwise. Thus $a[1:r]$ is obtained by drawing balls from the urn randomly, without replacement, under the condition that its number be not greater than r . This implies that every permutation of $\{1, \dots, r\}$ has probability $1/r!$ of actually occurring. Moreover, which ball with number greater (not greater) than r is drawn, is independent of the balls with number not greater (greater) than r drawn previously. Therefore, the permutation being formed in $a[1:r]$ is independent of the one being formed in $a[r+1:n]$ and vice versa.

This implies that the number of possibilities after the partition is $r!(n-r)!$, each with equal probability, so that the uncertainty is $\log_2 (r!(n-r)!)$. Introducing a quantity H , we find for the information yield of a partition, which equals the decrease of uncertainty,

$$nH = \log_2 (n!) - \log_2 (r!) - \log_2 ((n-r)!).$$

In quicksort, the bound y is selected in such a way that $r(y)$ has probability $1/n$ to become equal to $1, \dots, n$. The generalization that we will consider consists of introducing a different way of selecting the bound having a different probability distribution. We will regard r as a random variable, such that

$$\text{prob } \{r = r\} = f_r, \quad r = 1, \dots, n, \quad \sum_{r=1}^n f_r = 1.$$

This gives for the expected information yield of a partition

$$E(nH) = \log_2 (n!) - \sum_{r=1}^n f_r (\log_2 (r!) + \log_2 ((n-r)!)).$$

This equals, asymptotically for large n ,

$$E(nH) \sim n \log_2 (n) - \sum_{r=1}^n f_r (r \log_2 (r) + (n-r) \log_2 (n-r))$$

and

$$E(H) \sim - \sum_{r=1}^n f_r \left(\frac{r}{n} \log_2 \left(\frac{r}{n} \right) + \frac{n-r}{n} \log_2 \left(\frac{n-r}{n} \right) \right).$$

For large n , the sum may be replaced by an integral:

$$E(H) \sim - \int_0^1 g(x) (x \log_2 (x) - (1-x) \log_2 (1-x)) dx,$$

where $g(x)$ is the probability density function of $x = r/n$. In the sequel we shall confine our attention to symmetric distributions, that is, where $g(x) = g(1-x)$, and then we have:

$$E(H) \sim -2 \int_0^1 g(x) x \log_2 (x) dx. \quad (2)$$

In our derivation of the efficiency of a sorting algorithm, apparently a partition of a sequence of length n is the natural unit, irrespective of the way in which this partition has been effected. However, it has become usual (see Hoare [5]) to use the comparison between an element and the bound as the unit. We need n of these to complete a partition; so $E(H)$ may be interpreted as the average information yield of a comparison.

This result (2) holds asymptotically for large n for partitions in sequences of length n . However, completing the process of sorting requires partitions in sequences of any length not smaller than, say, 2. Suppose θ is the proportion of all comparisons required for subsequences of length $\leq pn$, where $(2/n) \leq p < 1$, and n is the length of the original sequence. θ attains its maximum θ_m when all partitions end exactly in the middle:

$$\theta_m = \frac{pn \log_2 (pn)}{n \log_2 (n)} = p \left(1 + \frac{\log_2 (p)}{\log_2 (n)} \right).$$

For any $p > 0$ we can choose an n large enough to satisfy $(2/n) \leq p$; hence $0 < \theta_m < p$ and also $0 < \theta < p$. This implies that for any fixed p ($0 < p < 1$) we can choose n so large that the proportion θ of all comparisons required for subsequences of length $\leq pn$ is smaller than p . Finally, we can make n sufficiently large to approach the asymptotic result (2) closely enough for sequences of length pn . This means, essentially, that, for $n \rightarrow \infty$, "almost all" comparisons are done in "large" sequences, and that expression (2) can be expected to apply to the total time taken by quicksort, not just to the first partitions.

Under our assumption of equally probable permutations of $1, \dots, n$, the total amount of information to be gained during sorting is $n \log_2 (n)$; so that we find for the average number of comparisons required for a generalization of quicksort:

$$T_n \sim n \log_2 (n) / E(H) = \alpha n \log_2 (n), \quad (3)$$

where

$$\alpha = \left(- \int_0^1 g(x)x \log_2(x) dx \right)^{-1}.$$

This formula, which was supplied, with a faulty proof, in van Emden [2], is due to F. E. J. Kruseman Aretz [6], who obtained it by a different method. It allows us to derive two results given by Hoare [5] as special cases.

The first result applies to the original version of quicksort. Here $g(x) = 1$ if $0 \leq x \leq 1$, $g(x) = 0$ otherwise. Substituted in (3), this yields $\alpha = 2 \ln(2) = 1.386$, which is Hoare's result.

The second result applies to the theoretical minimum of α . Consider the random variable \bar{x} that has $g(x)$ as probability density function. We shall derive the minimum of α where $g(x)$ is allowed to vary over all symmetric functions such that $\int_0^1 g(x) dx = 1$. One of the forms that Jensen's inequality might assume (see, for instance, [1]) is:

$$E(f(\bar{x})) \leq f(E(\bar{x})), \quad (4)$$

where \bar{x} is a random variable assuming nonnegative values with probability one, and f is a continuous and convex function.

If we choose $f(x) = -x \log_2(x)$, we find

$$\begin{aligned} E(-\bar{x} \log_2(\bar{x})) &= - \int_0^1 g(x)x \log_2(x) dx \\ &= 1/(2\alpha) \leq -E(\bar{x}) \log_2(E(\bar{x})). \end{aligned}$$

Because of the supposed symmetry of $g(x)$, $E(\bar{x}) = .5$; hence $\alpha \geq 1$, which is Hoare's result.

In Jensen's inequality, equality occurs if and only if \bar{x} assumes one value with probability 1. This value can only be .5 in the symmetric case; we may conclude that 1 is the lower bound for α , which is assumed if and only if $\text{prob}(\bar{x} = .5) = 1$, which corresponds to such a choice of the bound in quicksort that every partition ends exactly in the middle.

3. Two Ways of Increasing Efficiency

3.1. Bound as Median of a Random Sample. Hoare [5] suggests an improvement of his algorithm to be obtained by choosing as the bound the median of random sample of size $2k + 1$, and remarks that the resulting saving is very difficult to estimate. This modification has also been studied by Frazer and McKellar [3]. The formula (3) for the asymptotic number of comparisons is applicable to this case, in which it assumes a particularly simple form.

It is a well-known fact in order statistics that, whatever the distribution of the elements themselves, their ranks are uniformly distributed and the probability density function of the median of a sample of size $2k + 1$ from a uniform distribution is $g(x) = x^k(1-x)^k/B(k+1, k+1)$, where

B is the Beta function. To find α , we have to evaluate the integral

$$\begin{aligned} \int_0^1 g(x)x \ln x dx &= (1/B(k+1, k+1)) \int_0^1 x^{k+1} \\ &\quad \times \ln x(1-x)^k dx \quad (5) \\ &= (B(k+2, k+1)/B(k+1, k+1)) \\ &\quad \cdot (\psi(k+2) - \psi(2k+3)) \\ &= (\psi(k+2) - \psi(2k+3))/2, \end{aligned}$$

where ψ is the logarithmic derivative of the gamma function:

$$\psi(x) = \frac{d}{dx} \ln \Gamma(x).$$

From (3) and (5) we find that

$$\begin{aligned} \frac{1}{\alpha_k} &= \frac{\psi(2k+3) - \psi(k+2)}{\ln(2)} \\ &= \frac{\frac{1}{k+2} + \frac{1}{k+3} + \cdots + \frac{1}{2k+2}}{\ln(2)} \\ &= \frac{1 - \frac{1}{2} + \frac{1}{3} - \cdots + \frac{1}{2k+1} - \frac{1}{2k+2}}{\ln(2)}. \end{aligned}$$

This shows that $\alpha_k > 1$ and $\lim_{k \rightarrow \infty} \alpha_k = 1$, as indeed required by Jensen's inequality (4). See Table I.

TABLE I

k	0	1	2	3	4	5	6	∞
α_k	1.386	1.188	1.124	1.092	1.073	1.061	1.053	1.000

It may be useful to remark that α_k is independent of n . Therefore, for any size $2k + 1$ of a random sample, for sufficiently large n , the time required to find the median of the sample is an arbitrarily small proportion of the time required to complete a partition. Although the asymptotic properties of this method do not depend on the efficiency of the method used to find the median, the following note may be of interest.

Van Wijngaarden [8] proposed the following modification of quicksort for finding the median of a sequence. Instead of sorting each of the parts designated by a partition, only that one is sorted in which the median lies. In this way the median is found as the last of a sequence of nested intervals containing it. He showed that the number of comparisons required is, asymptotically for large n , equal to βn , where β does not depend on n . Kruseman Aretz [6] showed that $\beta = 2 + 2 \ln(2)$.

3.2. Bounding Interval. Quicksort starts each partition

by designating an element of $a[1:n]$ as the bound y . We found that efficiency may be improved by merely postponing the decision of what the bound is going to be. In fact, during the whole of the partition we only use an interval containing the bound; at any time any element of this interval could still be chosen without disturbing the partition obtained so far; hence the term "bounding interval" which we have chosen for this strategy.

To be specific, suppose that integers p and q satisfy the relations (1). In this situation there is no need to choose y as bound; but if, for instance,

$$xx = \max_i a[i], \quad 1 \leq i \leq p,$$

and

$$zz = \min_i a[i], \quad q \leq i \leq n,$$

then any element whose rank is not less than $r(xx)$ and not greater than $r(zz)$ might be chosen. To complete the partition, p must be increased and q must be decreased. It is possible to preserve the validity of (1) by, if necessary, interchanging elements, or by increasing xx , or by decreasing zz .

When at last $p + 1 = q$, the interval containing the bound has shrunk to the pair $\{xx, zz\}$, either of which could be chosen as bound. In actual fact neither is, because the necessity of a bound has vanished: the partition is complete already.

This particular way of effecting a partition implies a certain density function g . It depends on the supposition that the ranks represent a random permutation of the integers $1, \dots, n$. It is only useful to use g to compute efficiency if it holds not only for the initial partition but also for all successive ones, that is, when the partition leaves a random permutation in each of the parts. That this is the case may be seen as follows.

Let r' be the rank of the final value of zz . A particular element in the left half is replaced if and only if its rank is $\geq r'$, even though the tests deciding its replacement are based on values of xx and zz that are, in general, not yet equal to their final values. Thus the replacements are drawn randomly, without replacement, from the uniform distribution on $1, \dots, r' - 1$, and so are elements that remain in the left part.

4. Computing the Asymptotic Efficiency when Using a Bounding Interval

4.1. *The Succession of Intervals as a Random Walk on the Unit Half-square.* As shown in Figure 1, an interval may be represented by a point P on the unit square, where the coordinates are the rank (divided by n) of its endpoints xx and zz . Because $r(xx) < r(zz)$, P may only lie above the diagonal shown. The following description applies to the strategy embodied in procedure qsort. Picking the first x on the left that is greater than xx and the first z on the right that is smaller than zz corresponds to a uniformly distributed drawing of a point Q from the points of the

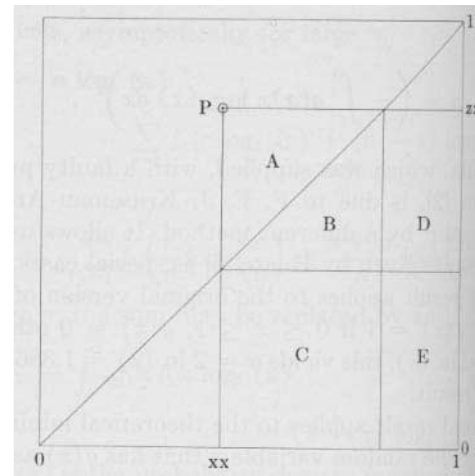


FIG. 1

unit square whose first (second) coordinate is greater (smaller) than that of xx (zz). These points form a rectangle of which P is the northwest corner; this rectangle may be referred to as the "shadow" of P .

The procedure qsort distinguishes whether Q is found in A, B, C, D, or E. If Q is found in A or B, both xx and zz are adjusted; if in C only zz , if in D only xx , and if in E neither xx nor zz are adjusted. Any of these adjustments is preceded by a reflection with respect to the diagonal if Q is found in B, C, D, or E. Under the assumption of a uniform distribution of Q on the shadow of P , the probability of each of these contingencies is proportional to its area. Hence we obtain the following succession rule:

If drawn from:	Probability	Q uniformly distributed on:
A or B	$(zz - xx)^2 u$	interior of A
C	$xx(zz - xx)u$	west side of A (7)
D	$(zz - xx)(1 - zz)u$	north side of A
E	$xx(1 - zz)u$	P

Here the factor $u = 1/(zz(1 - xx))$ ensures that these probabilities add up to 1.

This rule of succession defines a random walk on the upper half of the square with absorption on the diagonal. We may regard a partition as the following experiment:

- (1) A point is drawn at random, according to a uniform distribution, from the upper half of the square.
- (2) The successors according to the above rules are regarded as the stations of a random walk. (7)

Every random walk ends on the diagonal, and the probability density function resulting from this experiment is the function $g(x)$ needed to compute the asymptotic efficiency.

To approximate α we have set up the following discretized model of the experiment (7). The unit square is divided into a number of equal sized square cells, for every one of which we compute the probability mass (probability

that the point is in the cell under the condition that it started where it did) at each stage of the experiment. We pick a certain cell such that all cells above it or to the left of it are empty (contain zero mass).

In this model the succession rules (6) are interpreted as rules governing the distribution of mass from a certain cell over all others of the unit half square that are not above it or to the left of it. Suppose that according to these rules a certain proportion, say p , remains in the cell. Another application of the rules leaves p^2 , and so on. We require the result of infinitely many applications, and apparently this is achieved by computing the effect of a single application, multiplying all numbers so obtained by $1/(1 - p)$, and putting the mass of the selected cell equal to zero.

Let V be the set of cells not on the diagonal containing nonzero mass. As long as V is nonempty, it is possible to select an element of V that, according to the rules, may only give mass to other cells but may never receive any. Applications of the process described above to such a cell decreases the number of elements in V by 1. The computation continues until V is empty. Then all probability mass has diffused into the squares on the diagonal and their masses may be regarded as a discrete approximation of the density function g .

This computation was carried out with rules (6) corresponding to procedure qsort and it yielded $\alpha = 1.140$. Strictly, our analysis shows that the effectiveness of a partition is 18 percent greater than in quicksort, where $\alpha = 1.386$. Although qsort is very simple to program, a partition requires slightly more time than in the case of quicksort, because occasionally the bounding interval has to be adjusted. This probably explains an observed saving in computing time of about 15 percent.

Acknowledgments. The author is much indebted to Professor Dr. F. E. J. Kruseman Aretz of Philips Research Laboratories, Eindhoven, Netherlands, for valuable guidance in this matter; also to the referees, whose remarks caused the paper to be extensively revised.

RECEIVED OCTOBER, 1969; REVISED MARCH, 1970

REFERENCES

1. BECKENBACH, E. F., AND BELLMAN, R. *Inequalities*. Springer, New York, 1961.
2. VAN EMDEN, M. H. Iets quicker dan quicker. *Informatie 11* (1969), 30-32.
3. FRAZER, W. D., AND MCKELLAR, A. C. Samplesort: a sampling approach to minimal storage time sorting. In proc. of the Third Annual Princeton Conf. on Information Sciences and Systems, 1969, 276-280.
4. HOARE, C. A. R. Algorithm 64, Quicksort. *Comm. ACM* 4, 7 (July 1961), 321.
5. HOARE, C. A. R. Quicksort. *Comput. J.* 5 (1962), 10-15.
6. KRUSEMAN ARETZ, F. E. J. Private communication.
7. SHANNON, C. *The Mathematical Theory of Communication*. U. of Illinois Press, Urbana, Ill., 1963.
8. VAN WIJNGAARDEN, A. Private communication.
9. VAN EMDEN, M. H. Algorithm 402: Increasing the efficiency of quicksort. To appear in *Comm. ACM* 11 (Nov. 1970).

Algorithms

L. D. FOSDICK, Editor

ALGORITHM 392

SYSTEMS OF HYPERBOLIC P.D.E. [D3]

ROBERT R. SMITH AND DENNIS MCCALL (Recd. 7 Jan. 1969 and 17 June 1969)

US Naval Electronics Laboratory Center, San Diego, CA 92152

KEY WORDS AND PHRASES: hyperbolic p.d.e., characteristic, extrapolation, second order p.d.e., quasilinear p.d.e.

CR CATEGORIES: 5.17

DESCRIPTION:

CHARAC solves the initial value problem for the quasilinear hyperbolic system of equations

$$\begin{aligned} A_1 U_x + A_2 U_y + A_3 V_x + A_4 V_y &= H_1 \\ B_1 U_x + B_2 U_y + B_3 V_x + B_4 V_y &= H_2 \end{aligned} \quad (1)$$

in two independent variables X, Y and two unknown functions $U(X, Y), V(X, Y)$, where $A_i = A_i(X, Y, U, V), \dots, H_2 = H_2(X, Y, U, V)$. Specified data X_i, Y_i, U_i, V_i ($i=1, \dots, M$) given along a noncharacteristic curve Γ are used to find U and V at characteristic grid points in the entire characteristic cone associated with the initial curve. Values in the opposite characteristic cone can be computed by specifying the initial data points X_i, Y_i, U_i, V_i in the opposite order (X_1, Y_1, U_1, V_1 becomes X_M, Y_M, U_M, V_M , etc.).

If the system (1) is hyperbolic, it can be reduced to a normal form containing directional derivatives along two characteristic directions. The derivation of this normal form is given in Forsythe and Wasow [1, p. 38].

For (1) the normal form is

$$\begin{aligned} \left(\frac{dY}{dX} \right)_i &= \sigma_i, \\ R_i &= \left(\frac{\delta U}{\delta X} \right)_i + S_i \left(\frac{\delta V}{\delta X} \right)_i = T_i, \end{aligned} \quad i = 1, 2, \quad (2)$$

where $(\delta/\delta X)_i$ is the directional derivative along the characteristic with slope σ_i . Let $A = A_1 B_3 - A_3 B_1$, $C = A_2 B_4 - A_4 B_2$, $B = \frac{1}{2}(A_1 B_4 - A_4 B_1 - A_3 B_2 + A_2 B_3)$. Then the coefficients in (2) are given by

$$\sigma_i(X, Y, U, V) = \frac{B - (-1)^i (B^2 - AC)^{1/2}}{A}$$

$$R_i(X, Y, U, V) = A_1(B_1\sigma_i - B_2) - B_1(A_1\sigma_i - A_2),$$

$$S_i(X, Y, U, V) = A_3(B_1\sigma_i - B_2) - B_3(A_1\sigma_i - A_2),$$

$$T_i(X, Y, U, V) = H_1(B_1\sigma_i - B_2) - H_2(A_1\sigma_i - A_2).$$

The system (1) is called hyperbolic if $B^2 - AC > 0$ and if $R_1 S_2 - R_2 S_1 \neq 0$.

The subroutine CH VAR (XYUV, VAR) computes the values $\sigma_1, \sigma_2, R_1, R_2, S_1, S_2, T_1, T_2$ from A_i, B_i, H_i evaluated at the values X, Y, U, V given in the array XYUV. (The subroutine CH COEF giving A_i, B_i, H_i must be provided by the user, see Examples.) The computed values are returned in the array VAR of length 8. If σ_i, R_i, S_i, T_i are known to the user, he may provide his own routine CH VAR.

The system (2) is discretized by Massau's method, which is described in Forsythe and Wasow [1]. Given two adjacent points on the initial curve Γ , the nonparallel characteristics through the points intersect at a third point adjacent to the curve Γ . The values