

Purdue University

Purdue e-Pubs

Department of Computer Science Technical
Reports

Department of Computer Science

1976

Approximating the Number of Blocks Accessed in Data Base Organizations

S. B. Yao

Report Number:
76-184

Yao, S. B., "Approximating the Number of Blocks Accessed in Data Base Organizations" (1976).
Department of Computer Science Technical Reports. Paper 127.
<https://docs.lib.purdue.edu/cstech/127>

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries.
Please contact epubs@purdue.edu for additional information.

APPROXIMATING THE NUMBER OF BLOCKS ACCESSED
IN DATA BASE ORGANIZATIONS

S. B. Yao
Department of Computer Science
Purdue University
West Lafayette, Indiana 47907

CSD-TR 184
April 1976

(to appear in the Comm. ACM)



Approximating the Number of Blocks Accessed in Data Base Organizations

by

S. B. Yao

Key Words and Phrases: data base, inverted file organization, data base performance and measurement, information retrieval, query answering

CR Categories: 3.70, 3.72, 4.33

When data records are grouped into blocks in secondary storage, it is frequently necessary to estimate the number of blocks X_D accessed for a given query. In a recent paper [1], Cardenas gave the expression

$$X_D = m(1 - (1 - 1/m)^k), \quad (1)$$

assuming that there are n records divided into m blocks and that the k records satisfying the query are distributed uniformly among the m blocks. The derivation of the expression was left to the reader as an exercise.

Let us take a closer look at the expression. $(1 - 1/m)$ gives the probability that a particular block does not contain a particular record. If k records are selected independently, then the probability that a particular block not being "hit" is given by $(1 - 1/m)^k$. Hence $1 - (1 - 1/m)^k$ gives the probability that a particular block is "hit", and the expression follows.

The assumption that the k records are selected independently implies selection with replacement. Since a record may be selected more than once, the k records may not be distinct. This is not valid in the case of a query access which retrieves all k distinct records at one time. In fact, Rothnie and Lozano showed that the result of eq. (1) gives the lowerbound of the expected number of blocks accessed [2]. A more accurate analysis based on selection without replacement was given by Severance, but the precision problem makes the expression obtained computationally intractable (Appendix D in [3]). A similar approach by Siler results in a rather complicated recursive formula which can be computed (Appendix B in [4]). Using a different

approach, a simple closed form was obtained by Yao in a different context [5]. The resulting expression was used in several applications [5,6,7] to estimate the expected number of data blocks accessed. Comparing to the Cardenas' approximation, it is shown that this refinement is significant when the blocking factor n/m is small. For large blocking factors (e.g., $n/m \geq 10$), the error involved in Cardenas' approximation is practically negligible.

Theorem (Yao). Given n records grouped into m blocks ($1 < m \leq n$), each contains n/m records. If k records ($k \leq n - n/m$) are randomly selected from the n records, the expected number of blocks hit (blocks with at least one record selected) is given by

$$m \cdot \left(1 - \prod_{i=1}^k \frac{nd - i + 1}{n - i + 1} \right) \text{ where } d = 1 - \frac{1}{m} \quad (2)$$

Proof: Let X be a random variable representing the number of blocks hit and let I_j be a random variable where

$$I_j = \begin{cases} 1 & \text{when at least one record in the } j\text{-th} \\ & \text{block is selected,} \\ 0 & \text{otherwise.} \end{cases}$$

The j -th block has $p = n/m$ records and there are $n - p$ records not in the j -th block. The probability that no records are selected from the j -th block is

$$\frac{C_k^{n-p}}{C_k^n} \quad \text{or} \quad \frac{C_k^{nd}}{C_k^n} \quad \text{where } d = 1 - \frac{1}{m}$$

It follows that the expectation of I_j is

$$E(I_j) = 1 - \frac{C_k^{nd}}{C_k^n}$$

Hence the expected number of blocks hit is

$$\begin{aligned} E(X) &= \sum_{j=1}^m E(I_j) \\ &= m \cdot \left(1 - \frac{C_k^{nd}}{C_k^n} \right) \end{aligned}$$

Using the identity $C_y^x = \frac{x!}{y!(x-y)!}$, we have

$$\begin{aligned} E(X) &= m \cdot \left(1 - \frac{(nd)!(n-k)!}{n!(nd-k)!} \right) \\ &= m \cdot \left(1 - \prod_{i=1}^k \frac{nd-i+1}{n-i+1} \right) \end{aligned}$$

Q.E.D.

The following corollary of the theorem is obvious:

Corollary. If $k > n - n/m$ or $m = 1$, then all m blocks are hit.

It is interesting to note that eq. (1) is independent of the number of records n . We observe that the approximation of eq. (1) is good if $k \ll n$ and $m \ll n$. Intuitively, if n is large and each block contains many records, the selection with and without replacement makes little difference, especially when a small number of records are selected. In many actual applications it is true that $k \ll n$, otherwise sequential processing of the records can be considered. Using Siler's example, the results from Cardenas, Siler, and Yao are compared in Figure 1. The Siler's result for 300 records grouped into 20 blocks with a blocking factor of 15 matches Yao's result. The errors of Cardenas' expression are plotted in Figure 2.

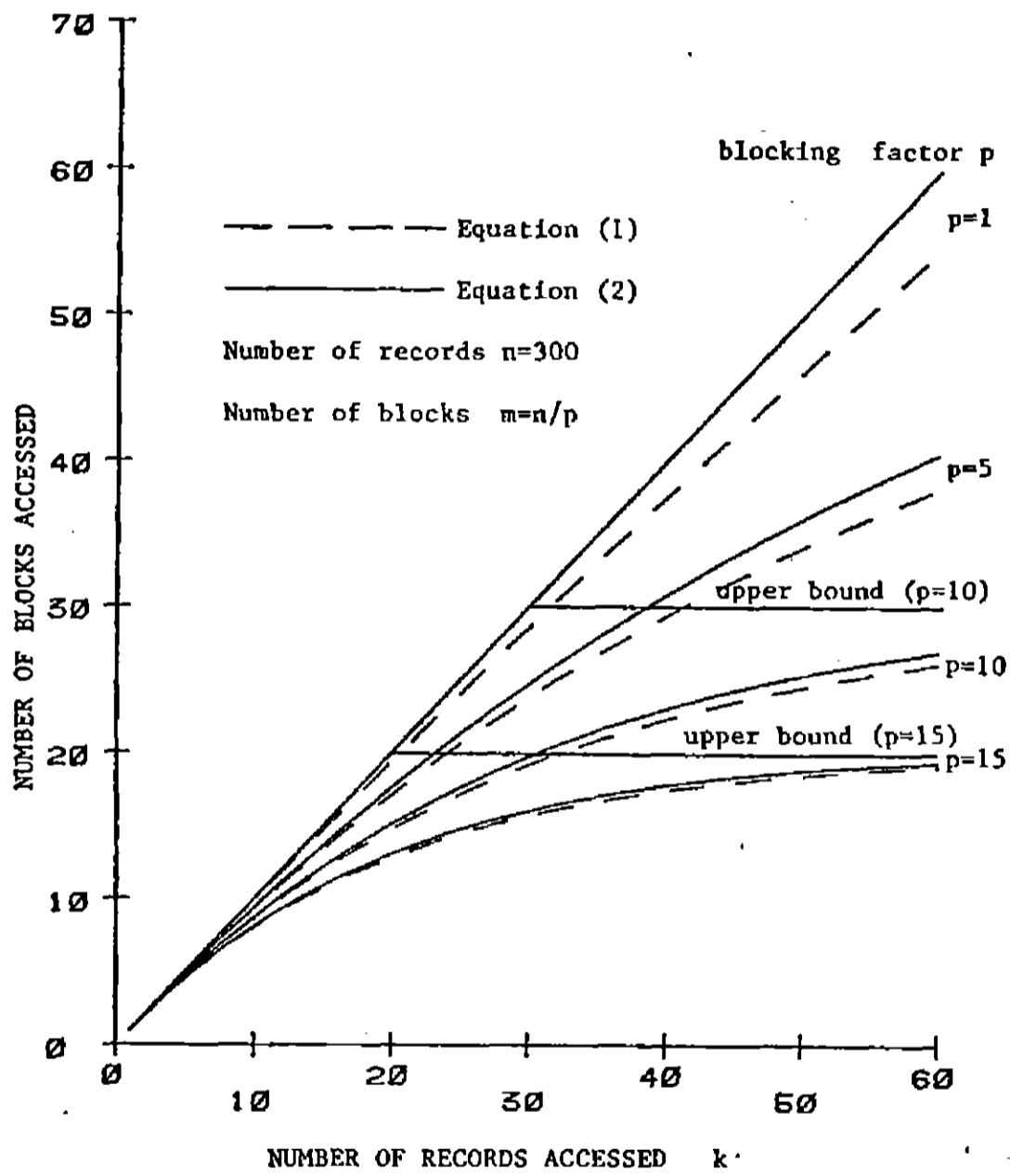


Figure 1. Comparison of expected block accesses.

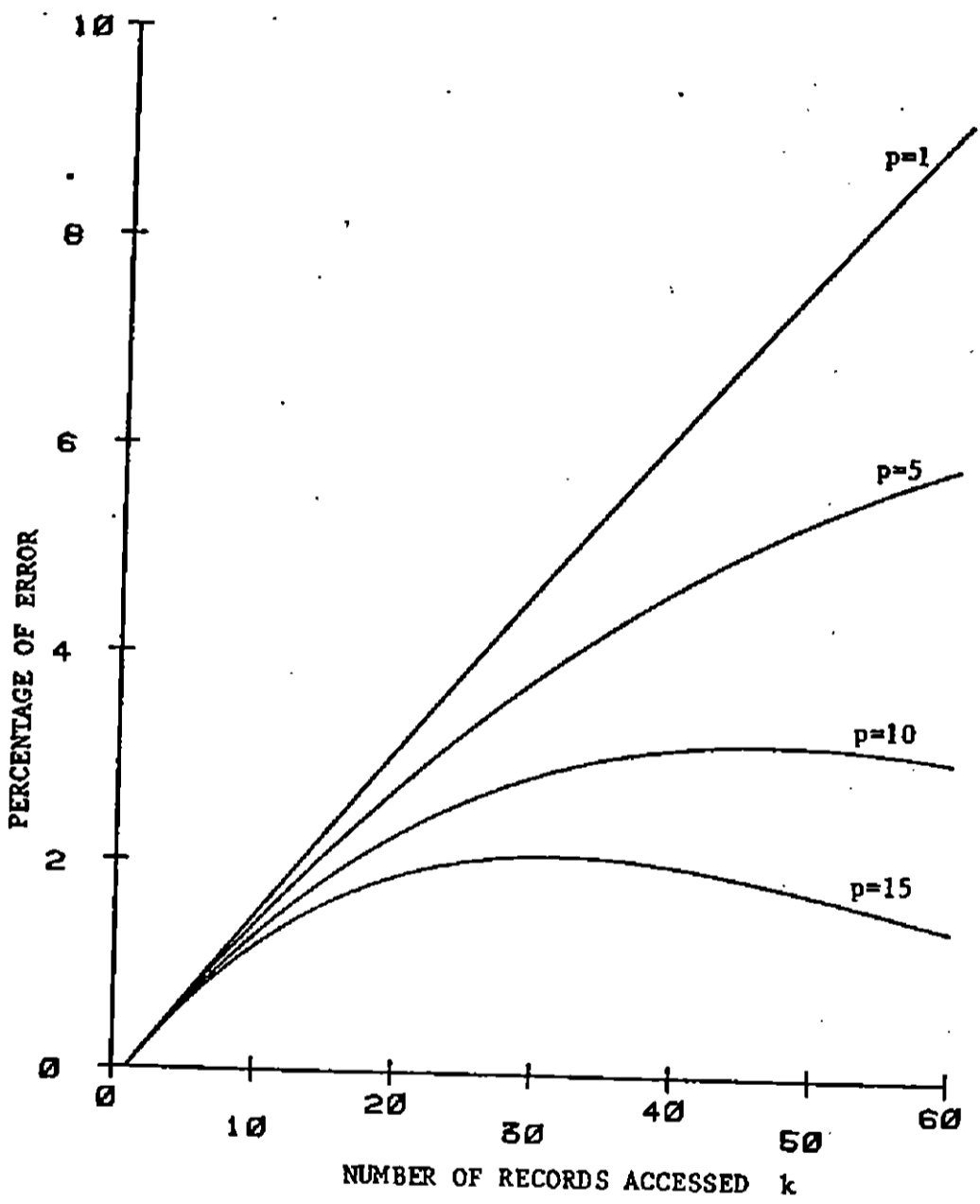


Figure 2. Errors of equation (1) approximation.

References

1. Cardenas, A. F. Analysis and performance of inverted data base structures. Comm. ACM 18, 5 (May 1975), 253-263.
2. Rothnie, J. B. and Lozano, T. Attribute based file organization in a paged memory environment. Comm. ACM 17, 2 (Feb. 1974) 63-69.
3. Severance, D. G. Some generalized modeling structures for use in design of file organizations, Ph.D. Dissertation, Univ. of Michigan, Ann Arbor, Michigan, 1972.
4. Siler, K. F. A stochastic evaluation model for data base organizations in data retrieval systems. Comm. ACM 19, 2 (Feb. 1976) 84-95.
5. Yao, S. B. Tree structures construction using key densities. Proc. ACM Nat. Conf. (Oct. 1975), 337-340.
6. Yao, S. B. Evaluation and optimization of file organizations through analytic modeling, Ph.D. Dissertation, Univ. of Michigan, Ann Arbor, Michigan, 1974.
7. Yao, S. B. A hierarchical access model for data base organizations. Technical Report TR-177, Computer Sciences, Purdue University, Lafayette, Indiana 47907, Feb. 1976.