

INFORMATION RETRIEVAL

Computational and Theoretical Aspects

H.S. HEAPS

INFORMATION RETRIEVAL

Computational and Theoretical Aspects

H. S. HEAPS

*Department of Computer Science
Concordia University
Montreal, Quebec
Canada*



ACADEMIC PRESS New York San Francisco London 1978
A Subsidiary of Harcourt Brace Jovanovich, Publishers

COPYRIGHT © 1978, BY ACADEMIC PRESS, INC.

ALL RIGHTS RESERVED.

NO PART OF THIS PUBLICATION MAY BE REPRODUCED OR TRANSMITTED IN ANY FORM OR BY ANY MEANS, ELECTRONIC OR MECHANICAL, INCLUDING PHOTOCOPY, RECORDING, OR ANY INFORMATION STORAGE AND RETRIEVAL SYSTEM, WITHOUT PERMISSION IN WRITING FROM THE PUBLISHER.

ACADEMIC PRESS, INC.

111 Fifth Avenue, New York, New York 10003

United Kingdom Edition published by

ACADEMIC PRESS, INC. (LONDON) LTD.

24/28 Oval Road, London NW1 7DX

Library of Congress Cataloging in Publication Data

Heaps, H S

Information retrieval, computational and theoretical aspects.

(Library and information science)

Includes bibliographical references.

1. Information storage and retrieval systems.

I. Title. II. Series.

Z699.H38 029.7 78-3338

ISBN 0-12-335750-0

PRINTED IN THE UNITED STATES OF AMERICA

Preface

The purpose of the present book is twofold. On the one hand, an attempt is made to introduce the student of computer science to some of the basic concepts of information retrieval and to describe the techniques required to develop suitable computer programs. On the other hand, an attempt is made to describe the general structure of the relevant computer programs so that basic design considerations may be understood by those not well versed in the details of computer science, such as librarians and information officers.

The material is organized with a view to create a textbook-like presentation rather than a comprehensive account of the state of the art. At the end of the chapters are problems intended to test the reader's understanding of the material and to lead the reader to consider further development of the basic principles.

The material in Chapters 1-4 and Chapters 6, 7, 10, and 11 has been found to be suitable for a one-term course for undergraduate computer science students who previously have covered the material in Chapter 5. Selected portions of Chapters 1-10 have been covered in a two-term course for students with a more general background. The material in Chapters 7-13 has also formed the basis of a course for graduate students.

It is believed desirable that the text should be reasonably self-sufficient so as to provide the non-computer-science student with a certain amount of

background material that is likely to be well known to the reader who has some previous education in computer science. This background material is contained at the beginning of Chapter 3 and in Chapter 5. Chapter 3 also serves to provide the student of computer science with illustrations chosen from the area of information retrieval.

Some mathematical maturity is required for a full understanding of Chapters 13 and 14. However, it is believed that a student not well-versed in mathematics should be able to appreciate the value of the techniques, and the significance of the results, without understanding the details of the analysis.

In Chapter 1 the reader is given a brief account of the motivation for development of the discipline of information retrieval from bibliographic data bases. A number of elementary general concepts are introduced in Chapter 2, and illustrations of the type of data base encountered in information retrieval are provided in Chapter 3. The portion of an information retrieval service that provides an interface between the user and the computer search is emphasized in Chapter 4 which deals with possible forms of question formulation. Thus Chapters 1-4 deal with aspects of information retrieval that are not directly concerned with the computational details and organization of the structure of the computer programs.

The structure of computer search programs for information retrieval is discussed in Chapter 6. This chapter refers to the data structures described in Chapter 5.

Chapters 7-10 are concerned with the vocabulary characteristics of document data bases and the manner by which a knowledge of vocabulary statistics may be used to advantage in the design of retrieval systems. Thus some general properties of bibliographic data bases are described in Chapter 7; the implications with respect to storage and transmission of information are discussed in a theoretical manner in Chapter 8. Some practical means of coding bibliographic data in order to economize in storage space are outlined in Chapter 9. Illustration of the application of the techniques to the design of a practical information retrieval system is given in Chapter 10.

The effectiveness of information retrieval from bibliographic data bases is dependent not only on the efficiency of the search programs but also on the extent to which the terms used to index the documents are good indicators of subject matter. This is discussed in Chapter 11. In Chapters 12 and 13 consideration is given to the problem of automatic selection of index terms and the related problems of automatic question modification and automatic document classification.

Contents

Preface	xi
---------	----

1

Introduction	1
--------------	---

1.1	Growth of Recorded Knowledge	1
1.2	The Discipline of Information Retrieval	4
1.3	Computer Learning and Adaptive Systems	8
1.4	Computer Identification of Meaning	9
1.5	Document Retrieval, Library Automation, and Privacy of Files	10

2

General Concepts	13
------------------	----

2.1	Document Data Bases and Selected Dissemination of Data	13
2.2	Coden Abbreviations	15
2.3	Library Data Bases	17
2.4	Numerical Data Bases	19
2.5	Management Information Systems	20
2.6	Keyword in Context and Keyword out of Context Indexes	20
2.7	Boolean Search	23

2.8.	Inverted Index and Double Dictionary	24
2.9	Precision and Recall	28
2.10	Thesauri	31
2.11	Terms and Vocabularies Associated with Attributes	33
2.12	Components of a Mechanized Information Retrieval System	35
2.13	Alphabetizing Conventions	35
2.14	Problems	37

3

Document Data Bases for Computer Search 39

3.1	Magnetic Tape and Disk Storage	39
3.2	Bit Codes for Data Storage	42
3.3	Blocks, Records, and Fields	47
3.4	Fixed and Variable Length Fields, Tags, and Directories	50
3.5	Example of Tagged Fields—METADEX Tape	53
3.6	Example of Fixed Length Tagged Fields— COMPENDEX Tape	57
3.7	Example of Fields with Noncharacter Tags—ERIC (AIM/ARM) Tape	58
3.8	Example of Tagged Fields and Subfields—SPIN Tape	59
3.9	Example of a Directory—CAIN Tape	63
3.10	Example of a Tagged Directory—MARC Tape	70
3.11	Preparation of Document Data Bases	76
3.12	Problems	79

4

Question Logic and Format 81

4.1	General Considerations	81
4.2	Truncation Specifications	83
4.3	Comparison and Termination Modes	84
4.4	Boolean Operations AND, OR, NOT WITH	85
4.5	The Ignore Specification	87
4.6	Adjacency and Precedence	88
4.7	Weighted Concepts	91
4.8	Defined Terms	93
4.9	Formal Description of the Question Syntax	95

4.10	Free Format for Question Formulation	97
4.11	User Specification of Output Format	98
4.12	Specification of Character Significance and Special Use	99
4.13	Problems	104

5

Data Structures for Storage and Retrieval 105

5.1	General Considerations	105
5.2	Sort Tree Structure	109
5.3	Dictionary Storage Using Character Tree	118
5.4	Table Structure to Allow Truncation Specification	120
5.5	Some Sorting Algorithms	126
5.6	Inverted File Structure	131
5.7	Scatter Storage	137
5.8	Stack Structures	142
5.9	Representation of Queues	145
5.10	List Storage Structure	147
5.11	Dynamic Storage	155
5.12	Problems	157

6

Structure of Search Programs 159

6.1	Sequential Search with Batched Questions	159
6.2	Single Nested OR Logic within AND Parameters	163
6.3	Question Processing through Logic Stack	167
6.4	Question Processing through Logic Tree	173
6.5	Question Processing through Inverted File	178
6.6	Problems	179

7

Vocabulary Characteristics of Document Data Bases 183

7.1	Dependence of Search Times on Vocabulary Characteristics	183
7.2	Vocabulary Frequencies	185
7.3	Distribution of Term Lengths	195

7.4	Frequency Distribution of Characters	199
7.5	Vocabulary Growth	206
7.6	Problems	208

8

Information Theory Considerations 209

8.1	Information Content of Textual Data	209
8.2	Information Content of Message Constraints	217
8.3	Information Gain of a Retrieval System	221
8.4	Huffman Codes for Compact Information Storage	224
8.5	Problems	227

9

Coding and Compression of Data Bases 229

9.1	Restricted Variable Length Term Codes	229
9.2	Hash Storage Based on Partial Keys	233
9.3	Coded Text Fragments	237
9.4	Partial Coding of Text	242
9.5	Term Compression by Abbreviation	243
9.6	Problems	245

10

Example of Design of a Document Retrieval System 247

10.1	Functional Description	247
10.2	Creation of Reformatted Tapes	249
10.3	Estimation of Data Base Statistics	251
10.4	Possible File Structure	252
10.5	File Update Procedure	255
10.6	Structure of Dictionary	256
10.7	Problems	260

11

Document Indexing and Term Associations 263

11.1	Document Representation through Index Terms	263
11.2	Assignment of Index Terms	264

11.3	Relative Frequencies of Terms in Text of Documents	269
11.4	Document Term and Term Connection Matrices	280
11.5	Term and Document Association Matrices	285
11.6	Information Retrieval through Stored Associations, Citation Indexing	290
11.7	Problems	291

12

Automatic Question Modification **293**

12.1	Weight and Response Vectors Related by Association Matrices	293
12.2	Automatic Question Modification through Association Feedback	298
12.3	Optimization of Retrieval Effectiveness	300
12.4	Further Discussion of the rms Search	306
12.5	Problems	308

13

Automatic Document Classification **309**

13.1	Document Classification by Categories	309
13.2	Attribute Analysis	310
13.3	Automatic Choice of Categories	315
13.4	Indexing Worth of Descriptors	316
13.5	A Measure of Classification or Retrieval Consistency	325
13.6	Problems	330

14

Concluding Remarks **333**

14.1	Limitations of Present Approach	333
14.2	Hardware Aspects	335
14.3	Theoretical Foundations	336

Subject Index	339
---------------	-----

1

Introduction

1.1 Growth of Recorded Knowledge

Information retrieval is a new discipline in the sense that many of its modern applications depend on concepts that have been formulated only during the past few decades. Nevertheless the subject has roots that extend back through many centuries.

The traditional library, as a collection of documents, led to the development of standard procedures for manual cataloguing, use of card indexes, bibliographies, and the circulation and ordering of books, journals, and reports. However, the traditional library was oriented more to the provision of documents than to the supply of information. This orientation is efficient provided that library users are interested primarily in well-defined subjects covered by a small number of books and journals. Yet at the present time there are many fields of study whose investigation requires information from a number of different disciplines, and for which requests for relevant information cannot be met by reference to a small, easily specified set of documents.

As remarked by C. P. Snow,¹ during all of human history until the present century the rate of social change has been sufficiently slow as to

¹ C. P. Snow, *The Two Cultures and the Scientific Revolution* (Cambridge: Cambridge University Press, 1959), p. 45.

pass unnoticed in one person's lifetime. This is no longer so. People are generally aware that social and technological changes are taking place very rapidly and are the result of discoveries that are understood by only a few specialists. Yet the consequences of many new discoveries affect the lives of entire populations to a degree that has never been the case previously. Thus more and more people are vitally interested in having fast access to more and more information.

The ability to maintain the rapid growth of technological and social changes requires that vast amounts of information be instantly available when required. The problem involved in meeting such a requirement may be emphasized by noting that it has been estimated that the number of scientific journals that existed in the year 1800, 1850, 1900, and 1966 was approximately 100, 1000, 10,000, and 100,000, respectively.² Holt and Schrank³ have indicated that between 1920 and 1960 the periodical literature in economics increased from 5000 to 40,000 articles per year. Likewise, the periodical literature in psychology increased from 30,000 to 90,000 articles per year. The annual number of papers published in mathematics during the period from 1868 to 1966 increased from about 800 to 13,000.^{4,5}

For a number of different disciplines the increases in the number of articles published in periodicals over the period between 1960 and 1970 have been estimated by Carter as shown in Table 1.1.⁶ Presumably the disagreement with Holt and Schrank's estimate for psychology is caused by a different definition of what constitutes a periodical article in psychology. The particular values of the individual figures are, however, less important than the implied growth rates of values estimated according to the same criteria.

The manner in which a new discovery may lead to an increase in the number of publications may be illustrated by noting the growth caused by the proposal of the laser structure in 1958.⁶ In 1960 there were approximately 20 papers on the subject of the ruby laser; in 1961 there were approximately 100 on the helium-neon laser; in 1962 there were 325 on the solid-state laser; in 1963 there were 700 papers on the GaAs laser,

² Proceedings of the Royal Institute of Great Britain, vol. 41, Part I, 1966.

³ C. C. Holt and W. E. Schrank, "Growth of the professional literature in economics and other fields, and some implications," *American Documentation* 19(1968):18-26.

⁴ K. O. May, "Quantitative Growth of the mathematical literature," *Science* 154 (1966):1672-1673.

⁵ K. O. May, "Growth and quality in the mathematical literature," *ISIS* 59(1968):363-371.

⁶ A. Neelameghan, "Theoretical Foundation for UDC: Its need and formulation," *Proceedings of the International Symposium*, Herceg Novi, Yugoslavia, June 28-July 1, 1971.

Table 1.1

Number of Articles Published in Periodicals During 1960 and 1970.

Subject	1960	1970
Mathematics	15,000	30,000
Physics	75,000	155,000
Civil engineering	15,000	15,000
Mechanical engineering	10,000	20,000
Electrical and electronic engineering	80,000	150,000
Aerospace engineering	35,000	75,000
Industrial engineering	15,000	15,000
Chemistry	150,000	260,000
Metallurgy	35,000	50,000
Biology	150,000	290,000
Geosciences	91,000	158,000
Agriculture	150,000	260,000
Medicine	220,000	390,000
Psychology	15,000	30,000
Other subjects	929,000	1,882,000
Totals	1,985,000	3,780,000

pulsed laser, and Q switching; in 1964 there were 1000 on the ion laser; and in 1965 there were 1200 papers on the N_2 - CO_2 high efficiency laser.

The term *information explosion* has been accepted very readily by workers in scientific fields who tend to be extremely conscious of the possibility of being unaware of work done previously, or being undertaken concurrently, by other scientists. There has even been a tendency, perhaps, to exaggerate the consequences of overlooking the work of others.⁷

The information needs of physicists and chemists engaged in research, administration, and teaching have been discussed in the literature.⁸ The roles of abstracting services in physics and summary papers in chemistry have been discussed respectively by Urquart⁹ and Bernal.¹⁰ A general discussion of the sources of information available to chemists and physi-

⁷ A. G. Oettinger, "An essay in information retrieval or the birth of a myth," *Information and Control* 8 (1965):64-79.

⁸ Survey of information needs of physicists and chemists. The report of a survey undertaken in 1963-4, in association with Professor B. H. Flowers, on behalf of the Advisory Council on Scientific Policy. *Journal of Documentation*, vol. 21, pp. 83-112, 1965.

⁹ D. J. Urquart, "Physics abstracting use and users," *Journal of Documentation* 21(1965):113-121.

¹⁰ J. D. Bernal, "Summary papers and summary journals in chemistry," *Journal of Documentation* 21(1965):122-127.

cists was given by Bottle.¹¹ The need for ready access to literature in the social sciences has been discussed by Guttman.¹²

1.2 The Discipline of Information Retrieval

Although recorded information usually is retrieved by means of stored data that represents documents, it is the emphasis on information relevant to a request, rather than direct specification of a document, that characterizes the modern subject of information retrieval. In addition to being concerned with the practical aspects of the design of operational computerized retrieval systems the subject of information retrieval includes aspects of the theory of measurement and definition, and of information content and relevance.

Some aspects of information retrieval may be compared to statistical communication theory as applied by electrical engineers and applied mathematicians during the past 30 years. The problem of locating relevant information from a body of widely dispersed knowledge is analogous to detection of the presence of a signal pulse in the presence of a noise background. Concepts such as the Wiener root-mean-square criterion, matched filters, feedback, and correlation detectors have their counterparts in the theory of information retrieval. It is perhaps rather curious that Norbert Wiener, who displayed great insight in the application of linear prediction techniques to problems in control theory and cybernetics, was very sceptical of the value of studies in information retrieval as he claimed that any information that he might require for study of his own subject could be obtained most easily by writing to any of the half dozen world experts in the field.

Wiener's point of view is interesting in that it serves to emphasize the difference between the environments of today's scholars and those of previous generations. Just as the number of people presently alive is a surprisingly large proportion of all people who have ever been born, so the amount of information in recorded form is many times larger than at any previous time. Moreover, not only is scientific, technological, and sociological information required by experts in these fields but also by nonspecialists who are unlikely to know the names of the leading authorities.

The rate of increase of available information has many philosophical

¹¹ R. T. Bottle, "A user's assessment of current awareness services," *Journal of Documentation* 21(1965):177-189.

¹² W. L. Guttman, "The literature of the social sciences and provision for research in them," *Journal of Documentation* 22(1966):186-194.

implications. The concept of man as a solitary traveller through time with some interaction from other human beings, and later the concept of man surrounded by a mechanistic universe, may now be replaced by the idea of man as an information receiver. Since information does not necessarily relate to physical quantities or to precisely measurable terms, it might be speculated that techniques developed for information retrieval and information evaluation eventually should be developed in a direction leading to further understanding of the process by which human beings associate ideas and gain understanding of scientific and humanistic concepts.

Spoken communication between humans is by sound waves that travel a distance of about 1000 feet in 1 sec. Communication between computers, or between computers and peripheral devices, is by electrical impulses that travel approximately 1000 feet in 1 μ sec. For a simple introduction to the principles of data transmission between computers, or through communication networks, the reader may refer to a paper of Kallenbach.¹³

Relays and switching devices within modern computers have response times of a few nanoseconds (1 nsec = 10^{-9} sec). Consequently, large amounts of information may be processed and transmitted by computers in a short time interval. Scientifically, well-informed, imaginative speculation on the possible consequences of utilizing such high transmission rates, and the ability of computers to absorb information at such speed, has led to many science-fiction writings such as, for example, those of the astronomer Hoyle.¹⁴⁻¹⁶ The degree to which instant communication has already affected contemporary society has been particularly emphasized from a popular viewpoint by McLuhan.¹⁷

Computerized information retrieval systems must be economical as well as feasible. In the same manner that economic considerations have led to more critical and hence more mathematically sophisticated design of engineering structures and chemical engineering processes, it is the economic considerations that are leading to a requirement for more precise mathematical formulation of the principles of information retrieval in order to ensure that the computers and computer accessible storage devices are used in an economic manner. Concurrently with the develop-

¹³ P. A. Kallenbach, "Introduction to data transmission for information retrieval," *Information Processing and Management* 11(1975):137-145.

¹⁴ F. Hoyle and G. Hoyle, *A for Andromeda* (Greenwich, Connecticut: Fawcett Publications Inc.).

¹⁵ F. Hoyle, *The Black Cloud* (London: Heinemann, 1957).

¹⁶ F. Hoyle, *October the First Is Too Late* (London: Heinemann, 1966).

¹⁷ M. McLuhan, *Understanding Media: The Extensions of Man* (New York: McGraw-Hill, 1964).

ment of new theory it is also important to study the behavior of operational retrieval systems in order to gain insight into the problems that arise from the point of view of the users of the system.

To some extent, the efficiency of a computerized information retrieval system is dependent on the computer hardware. Since computer hardware is continually being improved, and since computers become obsolete with the advent of more powerful and more economical ones, it is likely that operational retrieval systems will continue to be subject to constant revision and expansion. This fact presents less difficulty than might first be thought since it is the hardware and computer programs that change; the vast quantities of stored information may remain unchanged or be converted within the computer.

In analysis of complex or vaguely defined problems, the importance of good notation cannot be overemphasized. An efficient notation allows concepts to be stated clearly and allows data to be visualized as a whole, instead of as a number of isolated items. It also enables problems to be formulated in a precise manner with a true appreciation of any assumptions involved.

For formulation of information retrieval problems the concept of matrix transformation often allows statements to be made in a very compact form. The structure of written text, including the extent to which it is predictable, may be analyzed by probability theory.

Information theory, as developed in mathematical terms by Shannon,¹⁸ is concerned with information content in terms of the amount of information required for identification of symbols or words rather than in terms of the knowledge communicated by them. Randomness, or uncertainty or lack of knowledge, is measured in terms of entropy. The power and limitations of the theory have been discussed by a number of authors.¹⁹ Extension of Shannon's theory and computer analysis of written text to analyze information content, subject matter, and style, offers a number of problems to the student of information science. It also suggests the possibility of further cooperation between computer scientists and workers in other disciplines.

A study of information retrieval is necessarily concerned with optimization since it is desired to retrieve relevant items in the shortest possible time, or with minimum expense, or with maximum efficiency in regard to some estimate of relevance. The measure of relevance may be formulated

¹⁸ C. E. Shannon, "A mathematical theory of communication," *Bell System Technical Journal* 27(1948):379-423, 623-656.

¹⁹ P. L. Garvin (ed.), *Natural Language and the Computer* (New York: McGraw-Hill, 1963).

in mathematical terms and leads to considerations based on the mathematical theory of pattern recognition.

The subject of information retrieval is thus developing through application of matrix notation, probability theory, optimization techniques, pattern recognition, and systems analysis through which operations are represented by mathematical models that may be programmed on a computer.

The practical problems of information retrieval involve sufficient quantities of data that the introduction of any reasonable attempt to be systematic and independent of human intuition involves the handling of a large amount of data. The greater the degree of sophistication required, the greater the task of dealing with the information, and the greater the need for computer use since modern computers can store large quantities of data on computer accessible files and can examine the data very rapidly.

The systems analyst who designs a computer program for information retrieval or for some aspect of library automation first represents the entire operation by a mathematical model. The subsequent computer programming then proceeds in relation to this assumed model, which is described to the computer by means of a program written in some computer language. The resulting situation is that when the computerized system is in use, the librarians who use it tend to describe and evaluate it in terms of traditional library concepts while the computer analysts and programmers describe it using a different language and tend to think about the mathematical model rather than the operational system. It is clearly necessary for the two groups of people, those with traditional library backgrounds and those with computer science backgrounds, to have sufficient knowledge of each others' problems and the means of describing them. Otherwise, no effective communication may occur between the two groups. Ability to treat the problems from a wider viewpoint than generally acquired in the course of either a library science or computer science education is one of the skills expected of the information scientist.

The purpose of the present text is twofold. On the one hand an attempt is made to introduce the computer science student to some of the basic problems of information retrieval and to describe the techniques required to develop suitable computer programs. On the other hand an attempt is made to describe the general structure of the relevant computer programs so that basic design considerations may be understood by information officers and librarians not well versed in the details of computer science.

The ability of computers to perform complex arithmetic operations is more important for scientific computations than for information retrieval. The problems that arise in information retrieval are more concerned with identification, storage, rearrangement, and sorting of alphabetic data. It is