

THE COMPUTATION LABORATORY

Harvard University
Cambridge, Massachusetts

Scientific Report No. ISR-10
INFORMATION STORAGE AND RETRIEVAL
to
The National Science Foundation

Cambridge, Massachusetts
March 1966

Gerard Salton
Project Director



Copyright, 1965

By the President and Fellows of Harvard College

Use, reproduction, or publication, in whole or in part, is permitted
for any purpose of the United States Government.

DOCUMENT RETRIEVAL SYSTEMS - OPTIMIZATION AND EVALUATION

A thesis presented

by

Joseph John Rocchio, Jr.

to

The Division of Engineering and Applied Physics

in partial fulfillment of the requirements

for the degree of

Doctor of Philosophy

in the subject of

Applied Mathematics

Harvard University

Cambridge, Massachusetts

March 1966

PREFACE

This thesis presents a functional model for document retrieval systems. Based on this model, a search request optimization algorithm and an automatic classification algorithm applicable to efficient storage organization are derived. In addition some novel evaluation measures for document retrieval systems are proposed.

The writer is deeply indebted to Professor Gerard Salton for his guidance during the conception and development of this thesis, and for his confidence, encouragement, and continuing support.

Particular thanks are also due to Mr. Michael Lesk for his valuable assistance with respect to programming problems and with operational aspects of the SMART system, and to the other members of the SMART project at the Harvard Computation Laboratory. In addition the editorial assistance and moral support of Mary Rocchio, the author's wife, were of immeasurable value.

Finally the writer wishes to acknowledge the support of this research by the National Science Foundation under grant GN 360, and by the Bell Telephone Laboratories under contract with Harvard University.

TABLE OF CONTENTS

	<u>page</u>
Preface	v
List of Figures	xi
List of Tables	xiii
Synopsis	xv
 CHAPTER 1 INTRODUCTION	 1-1
1. The Document Retrieval Problem	1-1
2. A Functional Model	1-2
3. A Specific Model - The SMART System	1-6
A. Property Vector Indexing	1-6
B. Request Processing	1-7
C. Angular Distance Matching	1-7
D. Terminology	1-8
 CHAPTER 2 THE INDEXING FUNCTION	 2-1
1. Introduction	2-1
2. Manual Indexing	2-2
3. Automatic Indexing	2-2
A. The Statistical Approach	2-3
B. Semantic Techniques	2-4
C. Syntactic Techniques	2-6
4. The Structure of Index Representations	2-8
5. Optimizing the Index Transformation	2-12

TABLE OF CONTENTS (continued)

	<u>page</u>
CHAPTER 3 SEARCH REQUEST FORMULATION	3-1
1. Introduction	3-1
2. Request Formulation	3-3
3. Request Optimization	3-5
4. Relevance Feedback	3-11
5. The Case of No Relevant Documents	3-20
6. Experimental Results	3-21
A. Some Sample Search Requests	3-21
B. Average Results and Successive Iterations	3-32
C. Convergence	3-42
CHAPTER 4 THE QUERY-DOCUMENT MATCHING FUNCTION	4-1
1. The Comparison of Structured Operands	4-1
2. Storage Organization	4-7
3. Automatic Document Classification	4-12
4. Classification and Metric Searching	4-16
5. A Heuristic Classification Algorithm	4-19
A. Basic Concepts	4-19
B. Description of the Classification Algorithm	4-20
6. Experimental Results	4-36

TABLE OF CONTENTS (continued)

	<u>page</u>
CHAPTER 5 EVALUATION OF DOCUMENT RETRIEVAL SYSTEMS . . .	5-1
1. The General Problem	5-1
2. Evaluation Measures and the Collection of Statistics	5-2
A. The Idealized Experiment	5-2
B. Evaluation Statistics	5-6
C. Output Characterization	5-12
D. The Precision-Recall Tradeoff	5-15
3. The Use of Optimal Queries in Test Design . .	5-17
4. Cutoff-Independent Performance Indices . .	5-19
A. Derivation	5-19
B. Experimental Use	5-29
Appendix A - The SMART System	A-1

LIST OF FIGURES

<u>Figure</u>	<u>page</u>
2.1 Excerpt from the SMART Thesaurus	2-7
2.2 Alternative Property Space Index Representations . . .	2-11
3.1 Query Formulation - The User-System Interfaces . . .	3-2
3.2 Geometrical Representation of Relevance Feedback . . .	3-14
3.3 Relevant Document Images (\bar{r}_1 and \bar{r}_2) Which Have Only Slight Association	3-19
3.4 Hypothetical Document-Document Correlations Among an Assumed Set of Relevant Documents	3-19
3.5 A Typical Sample Search Request	3-26
3.6 Query Processing Using Relevance Feedback for Search Request "I-R Indexing"	3-29
3.7 Correlation Distributions for Search Request "I-R Indexing"	3-31
3.8 Query Processing Using Relevance Feedback for Search Request "Pattern Recognition"	3-33
3.9 Query Processing Using Relevance Feedback for Search Request "Analog-Digital"	3-34
3.10 Precision vs. Recall for Initial Queries and Queries Modified by Relevance Feedback (Averaged Over 24 Search Requests)	3-36
3.11 Comparison of Retrieval Results Using Original and Modified Query for Search Request "Random Numbers" . . .	3-37
3.12 Precision vs. Recall for Initial Queries and Two Iterations of Relevance Feedback	3-40
3.13 Retrieval Ranks of Relevant Documents for Query "Automata" Comparing Original Results with Those Obtained on Successive Iterations	3-41

LIST OF FIGURES (continued)

<u>Figure</u>	<u>page</u>
3.14 Retrieval Performance as a Function of Amount of Feedback for Query "I-R Indexing"	3-43
3.15 Retrieval Performance as a Function of Amount of Feedback for Query "M9 Natlang"	3-44
3.16 Document - Document Correlation Matricies	3-45
3.17 Retrieval Results for Two Queries Produced from Original Search Request "M9 Natlang"	3-46
4.1 Set Image Matching Operations	4-4
4.2 Retrieval Operations with Boolean Query Images	4-6
4.3 Total Number of Comparisions Required vs. the Number of Categories	4-22
4.4 Flowchart of Region Density Test	4-25
4.5 Graphical Illustration of the Density Test	4-26
4.6 Program for Specifying Cutoff Correlation	4-27
4.7 Progression of Categories and Correlation Distributions.	4-30
4.8 Classification Vectors	4-31
4.9 Flowchart of the Classification Algorithm	4-35
4.10 Cumulative Frequency Distribution of the Cutoff Correlation for Three Classifications	4-39
4.11 Query-Classification Vector Correlation Distribution for Request "Core Memory"	4-40
4.12 Evaluation of Two Level Searching for Test Query "Core Memory"	4-43
4.13 Average Evaluations for 24 Search Requests	4-44
5.1 Characterization of Retrieval Results	5-4

LIST OF TABLES

<u>Table</u>	<u>page</u>
2.1 Word Stem Frequency List of a Sample Document	2-10
3.1 Program Steps for Producing Relevance Feedback Modified Queries	3-23
3.2 Titles of Documents Relevant to Search Request "I-R Indexing"	3-30
4.1 Comparison Operations on Set Represented Operands . . .	4-3
4.2 Comparison Operations on Vector Represented Operands . .	4-7
4.3 Summary of the Steps of the Classification Algorithm . .	4-34
5.1 Retrieval results for 4 Equally Probable Query Types . .	5-11
5.2 Comparison of Precision and Recall Estimates	5-12

SYNOPSIS

A model for document retrieval systems consisting of the functional elements: indexing, search request formulation, and request-document matching is examined in this thesis. Analysis of the request formulation function leads to the definition of an optimality criterion and to a request optimization algorithm suitable for use in a system environment which allows iterative searching and real time user-system interaction. The optimality criterion has applicability to the evaluation of index languages, and generally to the design of evaluation tests for document retrieval systems. Investigation of the request-document matching leads to a novel automatic classification algorithm applicable to metric comparison measures, and useful for establishing an efficient storage organization. Finally the statistical basis for the evaluation of document retrieval systems is reviewed and some novel performance measures are proposed which are particularly suited to systems which induce a retrieval ordering on the members of the searched collection.

Chapter 1 is introductory in nature and attempts to define the area of discourse and its relation to the general field of information retrieval. A general model for document retrieval systems is introduced, and the basic functional elements of this model are briefly outlined. This material draws heavily on the work of Salton and the SMART automatic document retrieval project, as well as the general literature of the field.

The indexing component of the model is discussed in chapter 2 which is primarily descriptive in nature. The work of a number of researchers in this area is cited in describing the development and current trends of automatic content analysis. Particular emphasis is placed on concept vector indexing techniques which incorporate thesaurus type semantic associations. The SMART system generates document and query index images by such techniques and these were used in obtaining the experimental results presented in Chapters 3 and 4. The original contributions in chapter 2 are related to the techniques proposed for index image optimization.

A search request optimization algorithm analytically derived from an assumed optimality criterion is presented in Chapter 3. The optimization algorithm and the notion of request optimization by an iterative sequence of retrieval operations are original with the author. In addition the notion of testing index language devices by the use of optimal search requests is original. Experimental results illustrating the optimization process are presented. These were derived by a simulation which was coded and run on the IBM 7094 in conjunction with the SMART retrieval system. A search request formulation based on this optimization technique offers the promise of improving the performance of document retrieval systems, given the current state of development of computer time sharing and man-machine communication technology.

Chapter 4 presents an original automatic document classification algorithm, heuristically motivated by considerations of search efficiency and by the functional nature of query-document matching operations. This algorithm was coded in Fortran and run on

the IBM 7094 producing several classifications of a collection of 405 document index images (incorporated into the SMART system). Some experimental results illustrating the use of the classifications for improving search efficiency are presented. Document classification of the type described is novel, in that it is proposed as a direct adjunct to the query-document matching operation. Normally, automatic classification algorithms have been considered as replacements for manual classification or indexing.

The statistical basis for the evaluation of document retrieval systems is discussed in Chapter 5. Several of the topics considered are based on previous work which is cited by bibliographic reference. The organization and presentation as well as some of the conclusions drawn are original. In addition some novel performance statistics are derived which are particularly applicable to query-document matching operations possessing a high degree of discrimination such as the correlation measure of the assumed model. Each of the statistics derived is capable of describing overall system performance with a single parameter in contrast with several of the evaluation measures in current use.

DOCUMENT RETRIEVAL SYSTEMS - OPTIMIZATION AND EVALUATION