

C465

***THE SMART
RETRIEVAL SYSTEM***

***Experiments in Automatic
Document Processing***

Prentice-Hall
Series in Automatic Computation
George Forsythe, editor

- ANSELONE, *Collectively Compact Operator Approximation Theory and Applications to Integral Equations*
- ARBIB, *Theories of Abstract Automata*
- BATES AND DOUGLAS, *Programming Language/One*, 2nd ed.
- BAUMANN, FELICIANO, BAUER, AND SAMELSON, *Introduction to ALGOL*
- BLUMENTHAL, *Management Information Systems*
- BOBROW AND SCHWARTZ, editors, *Computers and the Policy-Making Community: Applications to International Relations*
- BOWLES, editor, *Computers in Humanistic Research*
- CRESS, DIRKSEN, AND GRAHAM, *FORTRAN IV with WATFOR and WATFIV*
- DANIEL, *The Approximate Minimization of Functionals*
- EVANS, WALLACE, AND SUTHERLAND, *Simulation Using Digital Computers*
- FIKE, *Computer Evaluation of Mathematical Functions*
- FORSYTHE AND MOLER, *Computer Solution of Linear Algebraic Systems*
- GAUTHIER AND PONTO, *Designing Systems Programs*
- GEAR, *Numerical Initial Value Problems in Ordinary Differential Equations*
- GORDON, *System Simulation*
- GREENSPAN, *Lectures on the Numerical Solution of Linear, Singular and Nonlinear Differential Equations*
- HARTMANIS AND STEARNS, *Algebraic Structure Theory of Sequential Machines*
- HEAPS, *An Introduction to Computer Languages*
- HULL, *Introduction to Computing*
- JOHNSON, *System Structure in Data, Programs, and Computers*
- LORIN, *Parallelism in Hardware and Software:*
An Introduction to Considerations in Real and Apparent Concurrency
- MARTIN, *Design of Real-Time Computer Systems*
- MARTIN, *Future Developments in Telecommunications*
- MARTIN, *Programming Real-Time Computer Systems*
- MARTIN, *Systems Analysis for Data Transmission*
- MARTIN, *Telecommunications and the Computer*
- MARTIN, *Teleprocessing Network Organization*
- MARTIN AND NORMAN, *The Computerized Society*
- MATHISON AND WALKER, *Computers and Telecommunications: Issues in Public Policy*
- MC KEEMAN, HORNING, AND WORTMAN, *A Compiler Generator*
- MINSKY, *Computation: Finite and Infinite Machines*
- MOORE, *Interval Analysis*
- PYLYSHYN, editor, *Perspectives on the Computer Revolution*
- RUSTIN, editor, *Computer Science*
- SALTON, *The SMART Retrieval System: Experiments in Automatic Document Processing*
- SAMMET, *Programming Languages: History and Fundamentals*
- SIMON AND SIKLOSSY, editors, *Representation and Meaning:*
Experiments with Information Processing Systems
- STERLING AND POLLACK, *Introduction to Statistical Data Processing*
- STROUD, *Approximate Calculation of Multiple Integrals*
- TAVISS, editor, *The Computer Impact*
- TRAUB, *Iterative Methods for the Solution of Equations*
- VARGA, *Matrix Iterative Analysis*
- VAZSONYI, *Problem Solving by Digital Computers with PL/I Programming*
- WILKINSON, *Rounding Errors in Algebraic Processes*

C465

GERARD SALTON, *Editor*

Professor of Computer Science
Cornell University

THE SMART
RETRIEVAL SYSTEM

***Experiments in Automatic
Document Processing***

PRENTICE-HALL, INC., Englewood Cliffs, New Jersey

BIBLIOTHEQUE DU CERIST

© 1971 by Prentice-Hall, Inc.
Englewood Cliffs, N.J.

All rights reserved. No part of this book
may be reproduced in any form or by any
means without permission in writing from
the publisher.

Current printing (last digit):

10 9 8 7 6 5 4 3 2 1

13-81425-3

Library of Congress Catalog Card No. 70-159122

Printed in the United States of America

PRENTICE-HALL INTERNATIONAL, INC., *London*
PRENTICE-HALL OF AUSTRALIA, PTY. LTD., *Sidney*
PRENTICE-HALL OF CANADA, LTD., *Toronto*
PRENTICE-HALL OF INDIA PRIVATE LIMITED, *New Delhi*
PRENTICE-HALL OF JAPAN, INC., *Tokyo*

PREFACE

The automatic SMART document retrieval system was designed at Harvard University between 1961 and 1964, and has been operating on IBM 7094 and 360 equipment both at Harvard and at Cornell University for several years. The system takes documents and search requests in the natural language, performs a fully automatic content analysis of the texts using one of several dozen programmed language analysis methods, matches analyzed documents with analyzed search requests, and retrieves for the user's attention those stored items believed to be most similar to the submitted queries.

Unlike the other computer-based retrieval systems, the SMART system does not rely on manually assigned keywords or index terms for the identification of documents and search requests, nor does it use primarily the frequency of occurrence of certain words or phrases included in the texts of documents. Instead, an attempt is made to go beyond simple word-matching procedures by using a variety of intellectual aids in the form of synonym dictionaries, hierarchical arrangements of subject identifiers, statistical and syntactic phrase generation methods, and the like, in order to obtain content identifications useful for the retrieval process.

By comparing the retrieval performance obtained with the various programmed procedures, the SMART system can be used as a unique experimental tool for the evaluation in a controlled laboratory environment of many fully automatic language analysis methods. In addition, the system

has been used to simulate a user-environment by making it possible for the user to participate in the search process. Specifically, the system utilizes feedback information supplied by the user during the search to construct improved search formulations, and to generate document representations reflecting the interests of the user population. By combining automatic text processing methods with interactive search and retrieval techniques, the SMART system may then lead to the design and implementation of modern information services of the type which may become current in operational environments some years hence.

Most of the documentation pertaining to the SMART system design and to the experimental results obtained with the system over the last few years is contained in a set of book-size scientific reports entitled "Information Storage and Retrieval," known in brief as "the ISR reports." The first ten of these (ISR-1 to ISR-10) were issued at Harvard between November 1961 and March 1966, and the last seven at Cornell. At the time of this writing, the most recent volume issued is ISR-17, dated September 1969.

The ISR reports covering the SMART system are not generally available in the open market; moreover, the information contained in the reports is difficult to assimilate, being dispersed over a large number of volumes including many thousands of pages. For this reason it has seemed advisable to collect in an organized manner, as a single book, the most important contributions contained in the earlier reports.

The present volume thus consists of updated versions of twenty-seven studies taken from the material contained in the ten most recent scientific reports (ISR-8 to ISR-17). Among the material covered are theoretical developments, including the derivation of system evaluation measures, language analysis techniques, document grouping techniques, and adaptive space transformation methods, as well as experimental studies relating to the evaluation of document analysis methods, interactive user feedback procedures, partial document searches based on clustered file organizations, and comparisons between the SMART system and more conventional operational information systems.

The present text is organized into eight major parts, entitled, respectively, The SMART System, Evaluation Viewpoint and Parameters, Language Analysis, Cluster Generation and Search, Basic Feedback Runs, Feedback Refinements, Document Space Transformation, and Operational Comparisons. Each part can be read independently of the remainder, and none of the material requires more than an elementary knowledge of mathematics or computer programming.

The text should be of greatest value as a reference volume for the professional practitioner interested in the design and operations of automatic information systems. It should also be useful as a text of readings in the area of automatic information retrieval for the more mature students enrolled in courses in applied mathematics, computer and information science, or library science. In this latter role, the studies contained in this volume might serve as a point of departure for term projects and for experimental work in modern information processing.

Part I of this volume, consisting of Chapters 1 and 2, covers the basic design of the SMART system. Chapter 1 by G. Salton deals with the background relating to the SMART project and contains a summary of evaluation results already obtained and a report of future plans. Chapter 2 by D. Williamson, R. Williamson, and

M. E. Lesk is a description of the implementation of the SMART system as it is presently operating at Cornell on an IBM 360/65.

The main systems evaluation parameters are covered in Part II, consisting of Chapters 3 to 5. Chapters 3 and 4 by J. J. Rocchio, Jr., contain a derivation of the normalized recall and precision evaluation measures and of the procedures used to obtain average retrieval results valid over many information searches. Chapter 5 by E. M. Keen is a detailed description of the design of a retrieval evaluation system.

The language analysis problems and relevant evaluation results are contained in Part III, consisting of Chapters 6 to 9. Chapter 6 by G. Salton and M. E. Lesk deals with the general problem of language analysis and dictionary construction useful in a content analysis system. Chapter 7 also by G. Salton and M. E. Lesk is a summarization of evaluation results obtained with SMART by processing document collections in the areas of aerodynamics, computer engineering, and documentation. A thorough analysis of the performance of user queries in the area of documentation is contained in Chapter 8 by E. M. Keen. The final chapter of Part III, number 9, by G. Salton covers an extension of the language analysis procedures originally implemented for English items to a collection of German language documents.

Part IV, consisting of Chapters 10 to 13, contains a description of automatic document classification procedures, and of the partial search methods based on document clusters. Chapter 10 by G. Salton is a basic description of cluster searching. Chapter 11 by R. T. Grauer and M. Messier covers the evaluation of a clustering process due to Rocchio. Additional, more efficient methods for automatic document classification are described in Chapter 12 by R. T. Dattola. Finally, Chapter 13 by S. Worona introduces an information search process based on the generation of request, rather than document, clusters.

The standard query transformation process based on user feedback is covered in Part V, consisting of Chapters 14 to 17. The main theoretical considerations relating to the construction of optimal user queries by relevance feedback methods are contained in Chapter 14 by J. J. Rocchio, Jr., and the principal feedback evaluation methods are described in Chapter 15 by G. Salton. A thorough analysis of feedback retrieval is included in Chapter 16 by E. Ide, and the last chapter, number 17, by C. Cirillo, Y. K. Chang, and J. Razon examines a number of novel feedback evaluation methods which circumvent some of the difficulties of the standard feedback evaluation methodology.

Certain refinements of the interactive search process are covered in Part VI of this volume, consisting of Chapters 18 to 22. The main distinctions between positive and negative feedback procedures are described in Chapter 18 by E. Ide and G. Salton; this chapter also contains a discussion of the relevance feedback process in the environment of a clustered information file. A query-splitting procedure, useful when nonhomogeneous sets of retrieved documents are identified as relevant during the search process is evaluated in Chapter 19 by A. Borodin, L. Kerr, and F. Lewis. A novel, experimental procedure for the implementation of negative feedback is contained in Chapter 20 by J. Kelly. Chapter 21 by J. S. Brown and P. D. Reilly deals with refined query modification methods in which the document and query terms receive individual treatment depending on the particular term characteristics. Finally Chapter 22 by D. Michelson, M. Amreich, G. Grissom, and E. Ide describes

procedures incorporating author information and bibliographic references into the feedback process.

The use of document, rather than query, transformations, made possible in most interactive retrieval environments is described in Part VII, consisting of Chapters 23 and 24. A procedure for document space transformation is introduced in Chapter 23 by S. R. Friedman, J. A. Maceyak and S. F. Weiss; a brief evaluation of the suggested procedure is also given. A modification of the transformation process covered in Chapter 23 is examined in detail in Chapter 24 by T. L. Brauen, and a thorough evaluation is presented both for the positive and the negative space modification processes.

The last part, number VIII, consists of Chapters 25 to 27, and covers various aspects connected with the operational evaluation of the SMART system. The use of a large variety of information displays of the kind presently available in an interactive retrieval environment with graphic or typewriter console equipment is described and evaluated in Chapter 25 by M. E. Lesk and G. Salton. Chapter 26 by M. E. Lesk and G. Salton covers an experiment designed to determine the importance of user relevance assessments for retrieval evaluation. Finally, in Chapter 27 by G. Salton a preliminary comparison is made between the retrieval effectiveness of the fully automatic SMART system, and the well-known MEDLARS system operating at the National Library of Medicine in Washington.

Readers wishing a quick overall view might restrict their attention to Parts I and III, and possibly to the operational problems described in Part VIII. The problems of content analysis are covered in detail in parts I, II, III and also VIII, whereas the interactive search procedures are contained in Parts IV, V, VI and VII. Readers interested in the problems of automatic language analysis should thus concentrate on the first set of chapters, while persons concerned with problems of file organization and file search should study the other half (Parts IV to VII).

The SMART system consisting of dozens of programmed routines and several hundred thousand machine instructions could not have been implemented without the help of the many programmers, analysts, and students who over the years have participated in the project, both at Harvard and at Cornell. Among those particularly instrumental who have extensively contributed to the design and programming phases are Dr. E. H. Sussenguth, Jr., now at IBM Corporation, Dr. J. J. Rocchio, Jr., now with International Computing Company, Dr. M. E. Lesk, now at Bell Telephone Laboratories; Mr. E. M. Keen, now at the College of Librarianship in Aberystwyth, Wales; and Mr. Robert E. Williamson, a graduate student at Cornell University. A number of other present or former students have performed important evaluation studies, including, in particular, T. Brauen, R. Dattola, and E. Ide.

To all these individuals, and to past and future users of the system, I am indebted for help and advice, and for extraordinary patience in bearing with the imperfections of an experimental system. I also wish to thank the McGraw-Hill Book Company for permission to include certain tabular material in Chapters 6 and 7. Finally, I am particularly grateful to the National Science Foundation whose continuous support over many years made it possible to design the system and to perform the research leading to the present implementation of the SMART system.

CONTENTS

PART I

THE SMART SYSTEM 1

1 THE SMART PROJECT—STATUS REPORT AND PLANS 3

G. Salton

- 1-1 Introduction 3
- 1-2 Experimental Results 4
- 1-3 Future Plans 6
 - 1-3-A. Text-Processing Experiments 6
 - 1-3-B. User Feedback and Document-Clustering Experiments 8
 - 1-3-C. Real-Time Operating System 10

2 THE CORNELL IMPLEMENTATION OF THE SMART SYSTEM 12

D. Williamson, R. Williamson, and M. E. Lesk

- 2-1 Introduction 12
- 2-2 Basic System Organization 13

- 2-2-A. Input of Printed Text 13
- 2-2-B. Document Clustering for Search Purposes 15
- 2-2-C. The Selection of Documents to be Searched 22
- 2-2-D. The Searching of the Document Groups 30
- 2-2-E. Search Evaluation 39
- 2-3 Access to the SMART System 47
- 2-4 Basic SMART System Flowchart 51

PART II

EVALUATION PARAMETERS 55

3 PERFORMANCE INDICES FOR DOCUMENT RETRIEVAL 57

J. J. Rocchio, Jr.

- 3-1 The Model 57
- 3-2 Evaluation Indices 59
- 3-3 Experimental Use 66

4 EVALUATION VIEWPOINTS IN DOCUMENT RETRIEVAL 68

J. J. Rocchio, Jr.

- 4-1 Introduction 68
- 4-2 Macro Evaluation 69
- 4-3 Micro Evaluation 70
- 4-4 Example 70
- 4-5 Conclusion 73

5 EVALUATION PARAMETERS 74

E. M. Keen

- 5-1 Purposes, Viewpoints, and Properties of Performance Measures 74
- 5-2 Measures for Ranking Systems 76
 - 5-2-A. Single Number Measures 77
 - 5-2-B. Varying Cut-off Performance Curves 78
 - 5-2-C. Comparison of Single Number and Curve Measures 80
- 5-3 The Construction of Average Precision Versus Recall Curves 84
 - 5-3-A. Averaging Techniques 84
 - 5-3-B. Cut-off Techniques 85
 - 5-3-C. Extrapolation Techniques for Request Generality Variations 91
 - 5-3-D. Extrapolation Techniques for Evaluation of Cluster Searching 96

- 5-4 Measures for Varying Relevance Evaluation 98
- 5-5 Measures for Varying Generality Comparisons 100
- 5-6 Techniques for Dissimilar System Comparisons and Operational Testing 103
- 5-7 The Comparison of Specific and General Requests and the Viewpoints of the High Precision and High Recall User 104
- 5-8 The Presentation of Data as Individual Request Merit 109

PART III

LANGUAGE ANALYSIS 113

6 INFORMATION ANALYSIS AND DICTIONARY CONSTRUCTION 115

G. Salton and M. E. Lesk

- 6-1 Introduction 115
- 6-2 Language Analysis 116
- 6-3 Dictionary Construction 119
 - 6-3-A. The Synonym Dictionary (Thesaurus) 119
 - 6-3-B. The Word-Stem Thesaurus and Suffix List 123
 - 6-3-C. The Phrase Dictionaries 127
 - 6-3-D. The Concept Hierarchy 130
- 6-4 Automatic Thesaurus Construction 132
 - 6-4-A. Fully Automatic Methods 133
 - 6-4-B. Semiautomatic Methods 135
- 6-5 Semiautomatic Hierarchy Formation 137

7 COMPUTER EVALUATION OF INDEXING AND TEXT PROCESSING 143

G. Salton and M. E. Lesk

- 7-1 Introduction 143
- 7-2 The SMART System 144
 - 7-2-A. Basic Organization 144
 - 7-2-B. Evaluation Process 146
 - 7-2-C. Significance Computations 150
- 7-3 Experimental Results 156
 - 7-3-A. Test Environment 156
 - 7-3-B. Document Length 158
 - 7-3-C. Matching Functions and Term Weights 159
 - 7-3-D. Language Normalization—The Suffix Process 166
 - 7-3-E. Synonym Recognition 167
 - 7-3-F. Phrase Recognition 168

- 7-3-G. Hierarchical Expansion 173
- 7-3-H. Manual Indexing 175
- 7-4 Concluding Comments 177

8 AN ANALYSIS OF THE DOCUMENTATION REQUESTS 181

E. M. Keen

- 8-1 Request Preparation 181
- 8-2 Characteristics of the Requests 182
 - 8-2-A. Length 182
 - 8-2-B. Important Request Words 183
 - 8-2-C. Multiple Need Requests 183
 - 8-2-D. Unclear Requests 183
 - 8-2-E. Difficult Requests 184
- 8-3 Relevance Decisions 185
- 8-4 Request Performance 186
 - 8-4-A. General Performance Analysis Methods 186
 - 8-4-B. Variation in Generality, Length, and Concept Frequency 187
 - 8-4-C. Comparison of Requests of the Two Preparers 191
 - 8-4-D. The Recognition of Important Request Words 194
- 8-5 Performance Effectiveness and Search Procedures 199

9 AUTOMATIC PROCESSING OF FOREIGN LANGUAGE DOCUMENTS 206

G. Salton

- 9-1 Introduction 206
- 9-2 The Evaluation of Language Analysis Methods 207
- 9-3 Multilingual Thesaurus 210
- 9-4 Foreign Language Retrieval Experiment 212
- 9-5 Failure Analysis 215
- 9-6 Conclusion 218

PART IV

CLUSTER GENERATION AND SEARCH 221

10 CLUSTER SEARCH STRATEGIES AND THE OPTIMIZATION OF RETRIEVAL EFFECTIVENESS 223

G. Salton

- 10-1 Introduction 223
- 10-2 Cluster Search Process 224
 - 10-2-A. Overall Process 224

- 10-2-B. Cluster Generation 226
- 10-2-C. Cluster Searching and Evaluation 234

11 AN EVALUATION OF ROCCHIO'S CLUSTERING ALGORITHM 243

R. T. Grauer and M. Messier

- 11-1 Introduction 243
- 11-2 A Description of Rocchio's Algorithm 244
- 11-3 Experimental Program 245
- 11-4 Evaluation System 246
 - 11-4-A. Tabulation of Results 247
 - 11-4-B. Detailed Analysis 251
 - 11-4-C. Conclusions and Remaining Questions 263

12 EXPERIMENTS WITH A FAST ALGORITHM FOR AUTOMATIC CLASSIFICATION 265

R. T. Dattola

- 12-1 Introduction 265
- 12-2 General Description 266
- 12-3 Implementation 268
 - 12-3-A. Initial Clusters 269
 - 12-3-B. Overlap 270
 - 12-3-C. Algorithm 272
- 12-4 Evaluation 275
 - 12-4-A. Evaluation Measures 275
 - 12-4-B. Internal Evaluation 281
 - 12-4-C. Initial Clusters 284
 - 12-4-D. Number of Clusters 287
 - 12-4-E. Overlap 289
 - 12-4-F. Cutoff 290
 - 12-4-G. Percent Loose Clustered 292
 - 12-4-H. External Evaluation 292
- 12-5 Conclusion 296

13 QUERY CLUSTERING IN A LARGE DOCUMENT SPACE 298

S. Worona

- 13-1 Introduction 298
- 13-2 Generating Clusters 299
- 13-3 Searching Clustered Collections 300
- 13-4 Parameters for Evaluating Cluster Searches 300
- 13-5 The Experiment 302
- 13-6 Results 304

PART V

BASIC FEEDBACK RUNS 311

14 RELEVANCE FEEDBACK IN INFORMATION RETRIEVAL 313*J. J. Rocchio, Jr.*

- 14-1 System Model 313
- 14-2 Request Formulation 314
- 14-3 Request Optimization 314
- 14-4 Relevance Feedback 316
- 14-5 Initial Experimental Results 318

15 RELEVANCE FEEDBACK AND THE OPTIMIZATION OF RETRIEVAL EFFECTIVENESS 324*G. Salton*

- 15-1 Relevance Feedback 324
- 15-2 Feedback Evaluation 326
- 15-3 Adaptive User-Controlled Multilevel Search 332

16 NEW EXPERIMENTS IN RELEVANCE FEEDBACK 337*E. Ide*

- 16-1 The Relevance Feedback Procedure 337
- 16-2 The Experimental Environment 339
- 16-3 Earlier Results in the Same Environment 340
- 16-4 Evaluation of Retrieval Performance 341
 - 16-4-A. The Feedback Effect in Evaluation 341
 - 16-4-B. Performance Measures 342
 - 16-4-C. Statistical Tests 344
- 16-5 Experimental Results 345
 - 16-5-A. Two Strategies Using Relevant Documents Only 345
 - 16-5-B. Varying the Amount of Feedback 346
 - 16-5-C. Strategies Using Nonrelevant Documents 348
- 16-6 Summary and Recommendations 353

17 EVALUATION OF FEEDBACK RETRIEVAL USING MODIFIED FREEZING, RESIDUAL COLLECTION, AND TEST AND CONTROL GROUPS 355*Y. K. Chang, C. Cirillo, and J. Razon*

- 17-1 Introduction 355
- 17-2 The Modified-Freezing Evaluation 356
- 17-3 Modified-Freezing Evaluation Results 358
- 17-4 The Residual Collection Method 359

- 17-5 Residual Collection Evaluation Results and Conclusions 362
- 17-6 The Test and Control Method 365
- 17-7 Test and Control Results and Evaluation 368
- 17-8 Conclusions 370

PART VI

FEEDBACK REFINEMENTS 371

18 INTERACTIVE SEARCH STRATEGIES AND DYNAMIC FILE ORGANIZATION IN INFORMATION RETRIEVAL 373

E. Ide and G. Salton

- 18-1 Retrieval System Performance 373
- 18-2 Request Space Modification 375
 - 18-2-A. Relevance Feedback 375
 - 18-2-B. Positive and Negative Strategies 376
 - 18-2-C. Selective Negative Feedback 383
- 18-3 Document Clustering 385
- 18-4 Document Space Modification 389
- 18-5 Conclusion 391

19 QUERY SPLITTING IN RELEVANCE FEEDBACK SYSTEMS 394

A. Borodin, L. Kerr, and F. Lewis

- 19-1 Introduction 394
- 19-2 The Query-Splitting Algorithm 395
- 19-3 Evaluation and Results 396
- 19-4 Conclusions and Suggestions for Further Research 401

20 NEGATIVE RESPONSE RELEVANCE FEEDBACK 403

J. Kelly

- 20-1 Introduction 403
- 20-2 Principal Algorithm 404
- 20-3 Experimental Method 405
- 20-4 Results 406
- 20-5 Conclusion 407

21 THE USE OF STATISTICAL SIGNIFICANCE IN RELEVANCE FEEDBACK 412

J. S. Brown and P. D. Reilly

- 21-1 Introduction 412
- 21-2 Query Construction 417

- 21-3 Conduct of the Experiment 418
- 21-4 Experimental Results 419
- 21-5 Conclusions and Recommendations 427

22 AN EXPERIMENT IN THE USE OF BIBLIOGRAPHIC DATA AS A SOURCE OF RELEVANCE FEEDBACK IN INFORMATION RETRIEVAL 430

D. Michelson, M. Amreich, G. Grissom, and E. Ide

- 22-1 Introduction 430
- 22-2 The Bibliographic Assumptions 431
- 22-3 The Problem and Method 431
- 22-4 Query Alteration in Feedback 434
- 22-5 Evaluation Results 435
- 22-6 Conclusions and Recommendations 441

PART VII

DOCUMENT SPACE TRANSFORMATIONS 445

23 A RELEVANCE FEEDBACK SYSTEM BASED ON DOCUMENT TRANSFORMATIONS 447

S. R. Friedman, J. A. Maceyak, and S. F. Weiss

- 23-1 The Problem 447
- 23-2 The Implementation 449
- 23-3 Experimental Results 451
- 23-4 Conclusions 455

24 DOCUMENT VECTOR MODIFICATION 456

T. L. Brauen

- 24-1 Introduction 456
- 24-2 Standard Document Vector Modification 458
- 24-3 Analysis of Standard Document Vector Modification 459
- 24-4 A First Proposal for Improved Document Space Modification 459
 - 24-4-A. The Method 459
 - 24-4-B. Testing Procedure 461
 - 24-4-C. Test Results 464
 - 24-4-D. Discussion 468
- 24-5 A Second Proposal for Improved Document Space Modification 474
 - 24-5-A. The Method 474
 - 24-5-B. Testing Procedure 476

24-5-C. Test Results	477
24-5-D. Discussion	481
24-6 Conclusions	484

PART VIII

OPERATIONAL COMPARISONS 485

25 INTERACTIVE SEARCH AND RETRIEVAL METHODS USING AUTOMATIC INFORMATION DISPLAYS 487

M. E. Lesk and G. Salton

25-1 Introduction	487
25-2 Fully Automatic Retrieval	489
25-3 User Interaction Through Presearch Methods	491
25-4 User Interaction Through Postsearch Methods	494
25-5 Evaluation Results and Discussion	495
25-5-A. Recall-Precision and Discussion	496
25-5-B. Overall Evaluation	501
25-6 Conclusion	504

26 RELEVANCE ASSESSMENTS AND RETRIEVAL SYSTEM EVALUATION 506

M. E. Lesk and G. Salton

26-1 Introduction	506
26-2 The Relevance Problem	507
26-3 The Experiment	510
26-4 Experimental Results	514
26-5 Judgment Consistency and Performance Measures	519
26-6 Machine-Search Effectiveness	521

27 A COMPARISON BETWEEN MANUAL AND AUTOMATIC INDEXING METHODS 528

G. Salton

27-1 Introduction	528
27-2 The Evaluation of Information Systems	529
27-3 The Test Design	532
27-3-A. The MEDLARS Evaluation Study	532
27-3-B. Design of the SMART Test	534
27-4 SMART-MEDLARS Comparison	537
27-5 Comparison of SMART Analysis Methods	541
27-6 Conclusions	545

INDEX 549