

INFORMATION RETRIEVAL

C. J. van RIJSBERGEN B.Sc., Ph.D., M.B.C.S.

Department of Computing Science

University of Glasgow

PREFACE TO THE SECOND EDITION

The major change in the second edition of this book is the addition of a new chapter on probabilistic retrieval. This chapter has been included because I think this is one of the most interesting and active areas of research in information retrieval. There are still many problems to be solved so I hope that this particular chapter will be of some help to those who want to advance the state of knowledge in this area. All the other chapters have been updated by including some of the more recent work on the topics covered. In preparing this new edition I have benefited from discussions with Bruce Croft, David Harper, Stephen Robertson and Karen Sparck Jones. I am grateful to the University of Cambridge Computer Laboratory for providing me with the facilities for carrying out the work. Finally, I am indebted to the Royal Society for supporting me on their Scientific Information Research Fellowship.

PREFACE TO THE FIRST EDITION

The material of this book is aimed at advanced undergraduate information (or computer) science students, postgraduate library science students, and research workers in the field of IR. Some of the chapters, particularly Chapter 6* , make *simple* use of a little advanced mathematics. However, the necessary mathematical tools can be easily mastered from numerous mathematical texts that now exist and, in any case, references have been given where the mathematics occur.

I had to face the problem of balancing clarity of exposition with density of references. I was tempted to give large numbers of references but was afraid they would have destroyed the continuity of the text. I have tried to steer a middle course and not compete with the *Annual Review of Information Science and Technology*.

* This is Chapter 7 in the second edition.

Normally one is encouraged to cite only works that have been published in some readily accessible form, such as a book or periodical. Unfortunately, much of the interesting work in IR is contained in technical reports and Ph.D. theses. For example, most the work done on the SMART system at Cornell is available only in reports. Luckily many of these are now available through the National Technical Information Service (U.S.) and University Microfilms (U.K.). I have not avoided using these sources although if the same material is accessible more readily in some other form I have given it preference.

I should like to acknowledge my considerable debt to many people and institutions that have helped me. Let me say first that they are responsible for many of the ideas in this book but that only I wish to be held responsible. My greatest debt is to Karen Sparck Jones who taught me to research information retrieval as an experimental science. Nick Jardine and Robin Sibson taught me about the theory of automatic classification. Cyril Cleverdon is responsible for forcing me to think about evaluation. Mike Keen helped by providing data. Gerry Salton has influenced my thinking about IR considerably, mainly through his published work. Ken Moody had the knack of bailing me out when the going was rough and encouraging me to continue experimenting. Juliet Gundry is responsible for making the text more readable and clear. Bruce Croft, who read the final draft, made many useful comments. Ness Barry takes all the credit for preparing the manuscript. Finally, I am grateful to the Office of Scientific and Technical Information for funding most of the early experimental work on which the book is based; to the Kings College Research Centre for providing me with an environment in which I could think, and to the Department of Information Science at Monash University for providing me with the facilities for writing.

C.J.v.R.

Contents

One: INTRODUCTION.....	1
The structure of the book.....	2
Outline.....	3
Information retrieval.....	3
An information retrieval system.....	4
IR in perspective.....	5
Effectiveness and efficiency.....	6
Bibliographic remarks.....	7
References.....	7
Two: AUTOMATIC TEXT ANALYSIS.....	10
Luhn's ideas.....	10
Generating document representatives - conflation.....	12
Indexing.....	13
Index term weighting.....	13
Probabilistic indexing.....	15
Discrimination and/or representation.....	17
Automatic keyword classification.....	17
Three: AUTOMATIC CLASSIFICATION.....	23
Measures of association.....	24
The cluster hypothesis.....	29
Single-link.....	35
The appropriateness of stratified hierarchic cluster methods.....	36
Single-link and the minimum spanning tree.....	38
Implication of classification methods.....	38
Conclusion.....	40
Bibliographic remarks.....	40
References.....	41
Four: FILE STRUCTURES.....	46
Introduction.....	46

Logical or physical organisation and data independence.....	46
A language for describing file structures.....	47
Basic terminology.....	47
Trees.....	58
Scatter storage or hash addressing.....	61
Bibliographic remarks.....	63
References.....	64
Five: SEARCH STRATEGIES.....	68
Introduction.....	68
Boolean search.....	68
Matching functions.....	69
Serial search.....	70
Cluster representatives.....	70
Cluster-based retrieval.....	73
Interactive search formulation.....	74
Feedback.....	75
Bibliographic remarks.....	77
References.....	77
Six: PROBABILISTIC RETRIEVAL.....	5
Introduction.....	5
Estimation or calculation of relevance.....	6
Basic probabilistic model*.....	7
Form of retrieval function.....	9
The index terms are not independent.....	10
Selecting the best dependence trees.....	12
Estimation of parameters.....	14
Recapitulation.....	16
The curse of dimensionality.....	17
Computational details.....	17
An alternative way of using the dependence tree (Association Hypothesis)	19
Discrimination power of an index term.....	20
Discrimination gain hypothesis.....	21

Bibliographic remarks.....	23
Seven: EVALUATION.....	5
Introduction.....	5
Relevance.....	6
Precision and recall, and others.....	7
Averaging techniques.....	9
Interpolation.....	10
Composite measures.....	12
The Swets model.....	12
Foundation.....	22
Presentation of experimental results.....	28
Significance tests.....	29
Bibliographic remarks.....	30
References.....	31
Eight: THE FUTURE.....	35
Future research.....	35
Automatic classification.....	35
File structures.....	36
Search strategies.....	36
Simulation.....	36
Evaluation.....	37
Content analysis.....	37
Future developments.....	38