

Textual and Visual Information Retrieval using Query Refinement and Pattern Analysis

S. G. Shaila · A. Vadivel

Textual and Visual Information Retrieval using Query Refinement and Pattern Analysis

 Springer

S. G. Shaila
Department of Computer Science
and Engineering
Dayananda Sagar University
Bangalore, India

A. Vadivel
Department of Computer Science
and Engineering
SRM University AP
Amaravati, Andhra Pradesh, India

ISBN 978-981-13-2558-8 ISBN 978-981-13-2559-5 (eBook)
<https://doi.org/10.1007/978-981-13-2559-5>

Library of Congress Control Number: 2018955166

© Springer Nature Singapore Pte Ltd. 2018

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Singapore Pte Ltd. The registered company address is: 152 Beach Road, #21-01/04 Gateway East, Singapore 189721, Singapore

This book is dedicated to my guide, who is the co-author of the book, my parents and my family. Their encouragement and understanding helped me to complete this work. I thank my husband and children, Vinu and Shamitha, for their love and support to bring this book as a reality.

Foreword

I am extremely delighted to write the foreword for *Textual and Visual Information Retrieval using Query Refinement and Pattern Analysis*. The authors, Dr. S. G. Shaila and Dr. A. Vadivel, have disseminated their knowledge on information retrieval through this book. The content of the book is a very good educational and valuable resource for researchers in the domain of information retrieval.

The book includes various chapters on deep Web crawler, event pattern retrieval, Thesaurus generation and query expansion, CBIR applications and indexing and encoding to cover the whole concept of information retrieval. Each chapter contains the theoretical information with experimental results and is intended for information retrieval researchers. The layout of each chapter includes a table of contents, introduction, material content and experimental results with analysis and interpretation.

This edition of the book reflects new guidelines that have evolved in information retrieval in terms of text- and content-based information retrieval schemes. The indexing and encoding mechanism of the low-level feature vector is also presented with results and analysis.

It is my hope and expectation that this book will provide an effective learning experience and referenced resources for all information retrieval researchers, leading to advanced research.

Bangalore, Karnataka, India

Dr. M. K. Banga
Chairman
Department of Computer
Science Engineering
Dean (Research)
School of Engineering
Dayananda Sagar University

Preface

Multimedia information retrieval from the distributed environment is an important research problem. It requires an architecture specification for handling various issues such as techniques to crawl information from WWW, user query prediction and refinement mechanisms, text and image feature extraction, indexing and encoding, similarity measure. In this book, research issues related to all the above-mentioned problems are discussed and suitable techniques are presented in various chapters.

Both the text and images are presented in web documents. In a comprehensive retrieval mechanism, text-based information retrieval (TBIR) plays an important role. In Chap. 1, the text-based retrieval is used for retrieving relevant documents from the Internet by using a suitable crawler with the capability to crawl deep and surface web. The functional dependency of core and allied fields in HTML FORM is identified for generating rules using SVM classifier. The presented crawler fetches a large number of documents while using the values in most preferable class. This architecture has a higher coverage rate and reduces fetching time.

In recent times, information classification is very important for text-based information retrieval. In Chap. 2, the classification based on events is presented and also the event Corpus is discussed, which is important for many real-time applications. Event patterns are identified and extracted at the sentence level using term features. The terms that trigger events along with the sentences are extracted from web documents. The sentence structures are analysed using POS tags. A hierarchal sentence classification model is presented by considering specific term features of the sentence, and the rules are derived along with fuzzy rules to get the importance of all term features of the sentences. The performance of the method is evaluated for 'Crime' and 'Die' and found that the performance of this approach is encouraging.

In general, the retrieval system depends on the user query to retrieve the web documents. The user-defined queries should have sufficient relevant terms, since the retrieval set depends on the queries. The query refinement through query expansion mechanism plays an important role. In Chap. 3, the N-gram Thesaurus construction mechanism for query expansion is presented. The HTML TAGs in web documents are considered and their syntactical context is understood. Based on the significance

of the TAGs in designing the web pages, suitable weight is assigned for TAGs. The term weight is calculated using corresponding TAG weight and frequency of the term. The terms along with the TAG information are updated into an inverted index. The N-grams are generated using the term and term weights in the document and updated as N-grams in the Thesaurus. During the query session, the term is expanded based on the content in the Thesaurus and suggested to the user. It is found that while the selected query is submitted to the retrieval system, the retrieval set consists of a large number of relevant documents.

In Chap. 4, the issues related to content-based image retrieval (CBIR) are presented. The chapter presents a histogram based on human colour visual perception by extracting the low-level features. For each pixel, the true colour and grey colour proportion are calculated using a suitable weight function. During histogram construction, the hue and intensity values are iteratively distributed to the neighbouring bins. The NBS distance is calculated between the reference bin and their adjacent bins. The NBS distance provides the overlapping proportion of the colour from the reference bin to its adjacent bins, and accordingly, the weight is updated. The distribution makes it possible to extract the background colour information effectively along with the foreground information. The low-level features of all the images in the database are extracted and stored in a feature database, and the relevant images are retrieved based on the rank. The Manhattan distance is used as a similarity measure, and the performance of the histogram is evaluated on Coral benchmark dataset.

In Chap. 5, the issues of indexing and encoding of low-level features and a similarity measure are presented. In CBIR system, the low-level features are stored along with the images and require a large number of storage space along with increased search and retrieval time. The search time increases linearly with the database size, which reduces the retrieval performance. The colour histograms of images are considered as low-level features. The bin values are analysed to understand their contribution representing image colour. The trivial bins are truncated and removed, and only important bins are considered to have histograms with lesser number of bins. The coding scheme used GR coding algorithm, and the quotient and remainder code parts are evaluated. Since there is variation between the number of bins in the query and database histogram, *BOSM* is used as a similarity measure. The performance of all the schemes is evaluated in an image retrieval system. The retrieval time, number of bits needed for histogram construction and precision of retrieval are evaluated using benchmark datasets, and the performance of the presented approach is encouraging.

Finally, as a whole, the book presents various important issues in information retrieval research filed and will be very much useful for the postgraduates and researchers working in information retrieval problems.

Bangalore, India
Amaravati, India

S. G. Shaila
A. Vadivel

Acknowledgements

First and foremost, we thank the Almighty for giving the wisdom, health, environment and people to complete this book.

We express our sincere gratitude to Dr. Hemachandra Sagar and Premachandra Sagar, Chancellor and Pro-Chancellor, Dayananda Sagar University, Bangalore; Dr. A. N. N. Murthy, Vice Chancellor, Dayananda Sagar University, Bangalore; Prof. Janardhan, Pro-Vice Chancellor, Dayananda Sagar University, Bangalore; Dr. Puttamadappa C., Registrar, Dayananda Sagar University, Bangalore; Dr. Srinivas A., Dean, School of Engineering, Dayananda Sagar University, Bangalore; Dr. M. K. Banga, Chairman, Department of CSE, and Dean Research, Dayananda Sagar University, Bangalore, for providing an opportunity and motivation to write this book.

We express our sincere gratitude to Dr. P. Sathyanarayanan, President, SRM University, Amaravati, AP; Prof. Jamshed Bharucha, Vice Chancellor, SRM University, Amaravati, AP; Prof. D. Narayana Rao, Pro-Vice Chancellor, SRM University, Amaravati, AP; Dr. D. Gunasekaran, Registrar, SRM University, Amaravati, AP, for providing an opportunity and motivation to write this book.

We would like to express our sincere thanks to our parents, spouse, children and faculty colleagues for their support, love and affection. Their inspiration gave us the strength and support to finish the book.

Dr. S. G. Shaila
Dr. A. Vadivel

Contents

1	Intelligent Rule-Based Deep Web Crawler	1
1.1	Introduction to Crawler	1
1.2	Reviews on Web Crawlers	2
1.3	Deep and Surface Web Crawler	4
1.4	Estimating the Core and Allied Fields	6
1.5	Classification of Most and Least Preferred Classes	8
1.6	Algorithm	9
1.7	Functional Block Diagram of Distributed Web Crawlers	10
1.8	Experimental Results	10
1.8.1	Rules for Real Estate Domain in http://www.99acres.com	12
1.8.2	Performance of Deep Web Crawler	14
1.9	Conclusion	17
	References	18
2	Information Classification and Organization Using Neuro-Fuzzy Model for Event Pattern Retrieval	21
2.1	Introduction	21
2.2	Reviews on Event Detection Techniques at Sentence Level	23
2.3	Schematic View of Presented Event Detection Through Pattern Analysis	26
2.4	Building Event Mention Sentence Repository from Inverted Index	26
2.5	Event Mention Sentence Classification	29
2.6	Refining and Extending the Rules Using Fuzzy Approach	33
2.6.1	Membership Function for Fuzzy Rules	34
2.6.2	Verification of Presented Fuzzy Rules Using Fuzzy Petri Nets (FPN)	37
2.7	Weights for Patterns Using Membership Function	40

2.8	Experimental Results	43
2.8.1	Performance Evaluation Using Controlled Dataset	43
2.8.2	Performance Evaluation for Uncontrolled Dataset Generated from WWW	48
2.8.3	Web Corpus Versus IBC Corpus	51
2.9	Conclusion and Future Works	53
	References	53
3	Constructing Thesaurus Using TAG Term Weight for Query Expansion in Information Retrieval Application	55
3.1	Introduction	55
3.2	Reviews on Query Expansion and Refinement Techniques	56
3.3	Architecture View of Query Expansion	59
3.4	TAG Term Weight (TTW)	60
3.5	N-Gram Thesaurus for Query Expansion	64
3.6	Algorithm	67
3.7	Experimental Results	69
3.7.1	Performance Evaluation of Query Reformulation and Expansion	69
3.7.2	Performance of TTW Approach	71
3.8	Conclusion	73
	References	73
4	Smooth Weighted Colour Histogram Using Human Visual Perception for Content-Based Image Retrieval Applications	77
4.1	Introduction	77
4.2	Reviews on Colour-Based Image Retrieval	79
4.3	Human Visual Perception Relation with HSV Colour Space	80
4.4	Distribution of Colour Information	81
4.5	Weight Distribution Based on NBS Distance	82
4.6	Algorithm	85
4.7	Experimental Results	86
4.8	Conclusion	90
	References	91
5	Cluster Indexing and GR Encoding with Similarity Measure for CBIR Applications	93
5.1	Introduction	93
5.2	Literature Review	94
5.2.1	Review on Indexing Schemes	94
5.2.2	Literature Review on Encoding Approaches	97
5.2.3	Literature Review on Similarity Metrics	98
5.3	Architectural View of Indexing and Encoding with Similarity Measure	99

5.4 Histogram Dimension-Based Indexing Scheme 99

5.5 Coding Using Golomb–Rice Scheme 102

5.6 Similarity Measure 103

 5.6.1 Algorithm—BOSM (H_q, H_k) 106

5.7 Experimental Results 107

 5.7.1 Experimental Results on Coding 107

 5.7.2 Retrieval Performance of BOSM on Various
 Datasets 108

 5.7.3 Evaluation of Retrieval Time and Bit Rate 109

 5.7.4 Comparative Performance Evaluation 110

5.8 Conclusion 116

References 116

Appendix 123

About the Authors

Dr. S. G. Shaila is an Associate Professor in the Department of Computer Science and Engineering in Dayananda Sagar University, Bangalore, Karnataka. She earned her Ph.D. in computer science from the National Institute of Technology, Tiruchirappalli, Tamil Nadu, for her thesis on multimedia information retrieval in distributed systems. She brings with her years of teaching and research experience. She worked for the DST project, Govt of India and Indo US based projects as Research Fellow. Her main areas of interest are information retrieval, image processing, cognitive science and pattern recognition. She published 20 international journals and conference proceedings.

Dr. A. Vadivel received his master's in science from the National Institute of Technology, Tiruchirappalli (NITT), Tamil Nadu, before completing a master's in Technology (M.Tech.) and Ph.D. at the Indian Institute of Technology (IIT), Kharagpur, India. He has 12 years of technical experience as a network and instrumentation engineer at the IIT Kharagpur and 12 years of teaching experience at Bharathidasan University and NITT. Currently, he is working as Associate Professor at SRM University, Amaravati, AP. He has published papers in more than 135 international journals and conference proceedings. His research areas are content-based image and video retrieval, multimedia information retrieval from distributed environments, medical image analysis, object tracking in motion video and cognitive science. He received the Young Scientist Award from the Department of Science and Technology, Government of India, in 2007; the Indo-US Research Fellow Award from the Indo-US Science and Technology Forum in 2008; and the Obama-Singh Knowledge Initiative Award in 2013.

Abbreviations

ACE	Automatic Content Extraction
ALNES	Active Long Negative Emotional Sentence
ALNIES	Active Long Negative Intensified Emotional Sentence
ALPES	Active Long Positive Emotional Sentence
ALPIES	Active Long Positive Intensified Emotional Sentence
ANFIS	Artificial neuro-fuzzy inference system
ANN	Artificial neural network
ASNES	Active Short Negative Emotional Sentence
ASNIES	Active Short Negative Intensified Emotional Sentence
ASPES	Active Short Positive Emotional Sentence
ASPIES	Active Short Positive Intensified Emotional Sentence
Bo1	Bose–Einstein statistics model
BoCo	Bose–Einstein statistics co-occurrence model
BOSM	Bin overlapped similarity measure
CART	Classification and regression tool
CBIR	Content-based image retrieval
CF	Core field
Co	Co-occurrence
CRF	Conditional random fields
DCT	Discrete cosine transform
DOM	Document object model
EMD	Earth mover’s distance
ET	Emotional triggered
FGC	Form Graph Clustering
FPN	Fuzzy Petri-Nets
GR	Golomb–Rice
HCPH	Human colour perception histogram
HiWE	Hidden Web Exposer
HQAOS	High-Qualified Active Objective Sentence
HQASS	High-Qualified Active Subjective Sentence

HQPOS	High-Qualified Passive Objective Sentence
HQPSS	High-Qualified Passive Subjective Sentence
HSV	Hue Saturation Value
HTML	Hyper-Text Markup Language
IBC	Iraq Body Count
InPS	Internal Property Set
IPS	Input Property Set
IR	Information retrieval
IRM	Integrated Region Matching
IWED	Integrated Web Event Detector
KDB	K-Dimensional B-tree
KLD	Kullback–Liebler divergence
KLDCo	Kullback–Liebler divergence co-occurrence model
LDA	<i>Latent Dirichlet allocation</i>
LP	Least Preferable
LVS	Label value set
MAP	Mean average precision
MCCM	Colour-based co-occurrence matrix scheme
ME	Mutually Exclusive
MEP	Minimum Executable Pattern
MHCPH	Modified human colour perception histogram
MP	Most Preferable
MRR	Mean reciprocal recall
MUC	Message Understanding Conference
NBS	National Bureau of Standards
NIST	National Institute of Standards and Technology
NLP	Natural language processing
NQAOS	Non-Qualified Active Objective Sentence
NQASS	Non-Qualified Active Subjective Sentence
NQPOS	Non-Qualified Passive Objective Sentence
NQPSS	Non-Qualified Passive Subjective Sentence
OPS	Output Property Set
PHOTO	Pyramid histogram of topics
PIW	Publicly Indexable Web
PLNES	Passive Long Negative Emotional Sentence
PLNIES	Passive Long Negative Intensified Emotional Sentence
PLPES	Passive Long Positive Emotional Sentence
PLPIES	Passive Long Positive Intensified Emotional Sentence
POS	Part of speech
PQ	Product quantization
PSNES	Passive Short Negative Emotional Sentence
PSNIES	Passive Short Negative Intensified Emotional Sentence
PSPES	Passive Short Positive Emotional Sentence
PSPIES	Passive Short Positive Intensified Emotional Sentence
QA	Question Answering

QAOS	Qualified Active Objective Sentence
QASS	Qualified Active Subjective Sentence
QPOS	Qualified Passive Objective Sentence
QPSS	Qualified Passive Subjective Sentence
RED	Retrospective new Event Detection
RS	Rule set
SCM	Sentence classification model
SOAP	Simple object access protocol
SQAOS	Semi-Qualified Active Objective Sentence
SQASS	Semi-Qualified Active Subjective Sentence
SQPOS	Semi-Qualified Passive Objective Sentence
SQPSS	Semi-Qualified Passive Subjective Sentence
SVM	Support vector machine
SWC	Surface web crawlers
SWR	Semantic Web Retrieval
TBIR	Text-based information retrieval
TDT	Topic Detection and Tracking
TS	Text Summarization
TSN	Term semantic network
TTW	TAG term weight
UMLS	Unified Medical Language System
URL	Uniform Resource Locator
ViDE	Vision-based data extraction
WaC	Web as Corpus
WSDL	Web Services Description Language
WWW	World Wide Web
XML	Extensible Markup Language

List of Figures

Fig. 1.1	Block diagram of deep web crawler	5
Fig. 1.2	Co-relation between core and allied fields.	7
Fig. 1.3	Classification of allied and core fields.	7
Fig. 1.4	Functional view of distributed web crawler.	10
Fig. 1.5	Sample core and allied fields	12
Fig. 1.6	a Rule for real estate web application. b Classification result.	13
Fig. 1.7	Average precision	15
Fig. 1.8	Coverage rate of the crawler. a Coverage rate for single fetch. b Coverage rate for periodic fetch	16
Fig. 1.9	Retrieval rate by crawlers	17
Fig. 2.1	Event features. a Relationship of Event features, b Sample Event features for crime-related document	22
Fig. 2.2	Presented approach framework	27
Fig. 2.3	Event Mention sentence repository	28
Fig. 2.4	Hierarchical Event Mention sentence classification	30
Fig. 2.5	Triangular membership function for <i>Subjective Active</i> class patterns.	36
Fig. 2.6	Triangular membership function for complete <i>Subjective</i> class patterns.	37
Fig. 2.7	Fuzzy Petri Net representation of the sentence classification model.	40
Fig. 2.8	Reachability graph	40
Fig. 2.9	Significant levels for patterns	41
Fig. 2.10	Weights derived for <i>Subjective</i> class patterns	42
Fig. 2.11	Event Corpus built from Event Instance patterns	43
Fig. 2.12	F1 score (%). a Various training data and b various Event Types	49
Fig. 3.1	Functional units of N-gram Thesaurus construction	59
Fig. 3.2	DOM tree for HTML TAGs	60
Fig. 3.3	Significant scale a <i><head></i> field and b <i><section></i>	62

Fig. 3.4	Weight for TAGs under <i><head></i> and <i><section></i>	63
Fig. 3.5	Weight distribution of HTML TAGs.	64
Fig. 3.6	Unigram Thesaurus	65
Fig. 3.7	Bigram Thesaurus	66
Fig. 3.8	Procedure to generate N-gram Thesaurus	67
Fig. 4.1	HSV colour model	81
Fig. 4.2	Distribution of true colour and grey colour components	82
Fig. 4.3	Construction of smooth weight distribution tree	83
Fig. 4.4	Smooth distribution of hue and intensity.	87
Fig. 4.5	Average precision	88
Fig. 4.6	Average recall.	88
Fig. 4.7	Average precision versus recall.	88
Fig. 4.8	Average <i>F</i> measure	89
Fig. 4.9	Sample retrieval set using MHCPH	89
Fig. 4.10	Sample retrieval set using HCPH	90
Fig. 5.1	Schematic diagram of the presented approach of indexing, encoding and distance similarity measure	100
Fig. 5.2	Sample histogram with empty cells (bins).	101
Fig. 5.3	Sample indexing structure.	101
Fig. 5.4	Values of <i>M</i> for various indexing level.	102
Fig. 5.5	View of database cluster.	104
Fig. 5.6	Working principle of BOSM.	105
Fig. 5.7	Retrieval performance of encoded and flat histogram a precision, b precision versus recall, c F1 score	108
Fig. 5.8	Performance of BOSM on MIT dataset_212 a precision, b precision versus recall	109
Fig. 5.9	Retrieval set MIT dataset_212 a clustered quotient code, b clustered combined, c flat <i>qcode</i> , d flat <i>ccode</i>	109
Fig. 5.10	Retrieval set from Caltech dataset_101	113
Fig. 5.11	Retrieval set from Caltech dataset_256	114
Fig. 5.12	Performance of similarity measure	114
Fig. 5.13	Performances of distance measures for Caltech dataset_101	115
Fig. 5.14	Performance on Caltech dataset_256.	115

List of Tables

Table 1.1	Estimated relationship between label and value for real estate web application.	11
Table 1.2	Combination of AF2 _j and AF1 for real estate web applications.	11
Table 1.3	Coverage Ratio.	17
Table 2.1	POS-tagged Event Mention sentences.	29
Table 2.2	POS tags used for sentence classification.	30
Table 2.3	Event Mention sentence classification in three levels using CART.	32
Table 2.4	Event Mention patterns.	32
Table 2.5	Sentential term features and POS tags	33
Table 2.6	Fuzzy rules of patterns for (a) <i>Subjective</i> class, (b) <i>Objective</i> class	35
Table 2.7	Sample Event Types and their Trigger terms	44
Table 2.8	IBC Corpus Statistics for ‘Die’ Event Type	44
Table 2.9	Classification accuracy of human annotation and ANFIS for Event Type ‘Die’	46
Table 2.10	Performance of ANFIS using tenfold cross-validation for Event Type ‘Die’	47
Table 2.11	Performance evaluations (%) for IBC dataset using k-fold cross-validation.	47
Table 2.12	Performance measure (%) of presented approach with other approaches for Event Type ‘Die’	48
Table 2.13	Web Corpus Statistics from www.trutv.com/library/crime	50
Table 2.14	F1 measure of the presented approach with other approaches for various Event Types in Web Corpus.	50
Table 2.15	F1 measure for various combinations of training/test data for the ‘Die’ Event	51
Table 2.16	Accuracy for the Instances for various Event Types.	53

Table 3.1	Information on benchmark dataset	69
Table 3.2	Clueweb09B: expanded queries	70
Table 3.3	Query expansion and user analysis.	70
Table 3.4	Baseline performance on Clueweb09B, WT10g and GOV2 datasets.	71
Table 3.5	Comparative performance of TTW.	71
Table 3.6	Performance comparison of various approaches for various datasets against baselines	72
Table 3.7	Gain improvement	73
Table 3.8	Difference in the improvement gain of comparative approaches with TTW approach.	73
Table 4.1	NBS distance table.	83
Table 5.1	Sample value of multiplication factor.	102
Table 5.2	Encoded feature (histogram).	103
Table 5.3	Sample distance value	104
Table 5.4	<i>BOSM</i> between histograms.	105
Table 5.5	Details of benchmark datasets.	107
Table 5.6	Retrieval time.	110
Table 5.7	Bit ratio for MIT dataset.	110
Table 5.8	Performance comparison on Caltech dataset_101	111
Table 5.9	Performance comparison on Caltech dataset_256	112
Table 5.10	Precision on Caltech dataset_101	116
Table 5.11	Precision on Caltech dataset_256	116