

Springer Theses

Recognizing Outstanding Ph.D. Research

Zechao Li

Understanding- Oriented Multimedia Content Analysis

 Springer

Springer Theses

Recognizing Outstanding Ph.D. Research

Aims and Scope

The series “Springer Theses” brings together a selection of the very best Ph.D. theses from around the world and across the physical sciences. Nominated and endorsed by two recognized specialists, each published volume has been selected for its scientific excellence and the high impact of its contents for the pertinent field of research. For greater accessibility to non-specialists, the published versions include an extended introduction, as well as a foreword by the student’s supervisor explaining the special relevance of the work for the field. As a whole, the series will provide a valuable resource both for newcomers to the research fields described, and for other scientists seeking detailed background information on special questions. Finally, it provides an accredited documentation of the valuable contributions made by today’s younger generation of scientists.

Theses are accepted into the series by invited nomination only and must fulfill all of the following criteria

- They must be written in good English.
- The topic should fall within the confines of Chemistry, Physics, Earth Sciences, Engineering and related interdisciplinary fields such as Materials, Nanoscience, Chemical Engineering, Complex Systems and Biophysics.
- The work reported in the thesis must represent a significant scientific advance.
- If the thesis includes previously published material, permission to reproduce this must be gained from the respective copyright holder.
- They must have been examined and passed during the 12 months prior to nomination.
- Each thesis should include a foreword by the supervisor outlining the significance of its content.
- The theses should have a clearly defined structure including an introduction accessible to scientists not expert in that particular field.

More information about this series at <http://www.springer.com/series/8790>

Zechao Li

Understanding-Oriented Multimedia Content Analysis

Doctoral Thesis accepted by
University of Chinese Academy of Sciences, Beijing, China

 Springer

Author
Dr. Zechao Li
Nanjing University of Science
and Technology
Nanjing
China

Supervisor
Prof. Hanqing Lu
National Laboratory of Pattern Recognition,
Institute of Automation
Chinese Academy of Sciences
Haidian District, Beijing
China

ISSN 2190-5053

Springer Theses

ISBN 978-981-10-3688-0

DOI 10.1007/978-981-10-3689-7

ISSN 2190-5061 (electronic)

ISBN 978-981-10-3689-7 (eBook)

Library of Congress Control Number: 2017939539

© Springer Nature Singapore Pte Ltd. 2017

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.

Printed on acid-free paper

This Springer imprint is published by Springer Nature

The registered company is Springer Nature Singapore Pte Ltd.

The registered company address is: 152 Beach Road, #21-01/04 Gateway East, Singapore 189721, Singapore

This work is dedicated to my parents, and all my friends. Their support and encouragement keep me forward.

Supervisor's Foreword

Multimedia content analysis has attracted extensive attention in the multimedia and social media research communities. Its goal is to reveal the semantic information intelligently. Zechaos' Ph.D. work focuses on understanding-oriented multimedia content analysis from the low-level visual representation to the high-level semantic understanding. As a key member of my group, he made a number of significant contributions in his research work. He investigated advanced multimedia content analysis approaches and proposed understanding-oriented multimedia content analysis approaches, including data representation (feature selection and feature extraction), tag recommendation, and multimedia news services. He directly integrated the visual understanding and learning models into a unified framework. The visual understanding guides the model learning while the learned models improve the visual understanding. The inspiring idea of understanding-oriented multimedia content analysis has been recognized as opening up possibilities to challenging multimedia content and context understanding. The proposed structured subspace learning framework has been successfully generalized to social image understanding, (semi-) supervised classification and clustering. His work has brought in new thoughts and disruptive models in understanding multimedia data. I believe that this book will benefit researchers and students conducting research on multimedia computing and social multimedia analysis.

Beijing, China
January 2017

Prof. Hanqing Lu

Preface

The amount of today's multimedia contents explosively grows due to the popularization and rapid growth of digital mobile devices and social media tools. To efficiently analyze and understand the multimedia content is still a challenging task. Over the past decade, many advanced methods have been proposed in the literature, including a few books on this topic. However, there is no book offering a systematic introduction to multimedia content analysis towards an understanding-oriented approach. Therefore, this book will focus on a novel "understanding" framework for multimedia content interpretation. This book offers a systematic introduction to multimedia content analysis towards an understanding-oriented approach. It integrates the visual understanding and learning models into a unified framework, within which the visual understanding guides the model learning while the learned models improve the visual understanding. More specifically, the book presents multimedia content representations and analysis including feature selection, feature extraction, image tagging, user-oriented tag recommendation, and understanding-oriented multimedia applications. By providing the fundamental technologies and the state-of-the-art methods, this book will be of interest to graduate students and researchers working in the field computer vision and machine learning.

Chapter 1 introduces the background, challenges, and progresses of understanding-oriented multimedia content analysis. Chapters 2 and 3 introduce some works of understanding-oriented data representation. The personalized tag recommendation work is detailed in Chap. 4, followed by understanding-oriented multimedia news services in Chaps. 5 and 6. Chapter 7 concludes the book by summarizing the major points and identifying the future works.

Nanjing, China
January 2017

Zechao Li

Parts of this book have been published in the following articles:

- Zechao Li, Jing Liu, Jinhui Tang, Hanqing Lu. Robust Structured Subspace Learning for Data Representation. *IEEE Trans. on Pattern Analysis and Machine Intelligence* 37(10): 2085–2098, 2015.
- Zechao Li, Jinhui Tang. Unsupervised Feature Selection via Nonnegative Spectral Analysis and Redundancy Control. *IEEE Trans. on Image Processing* 24(12): 5343–5355, 2015.
- Zechao Li, Jing Liu, Yi Yang, Xiaofang Zhou, Hanqing Lu. Clustering-Guided Sparse Structural Learning for Unsupervised Feature Selection. *IEEE Trans. Knowledge and Data Engineering* 26(9): 2138–2150, 2014.
- Jing Liu, Zechao Li, Jinhui Tang, Yu Jiang and Hanqing Lu. Personalized Geo-Specific Tag Recommendation for Photos on Social Websites. *IEEE Trans. on Multimedia* 16(3): 588–600, 2014.
- Zechao Li, Jinhui Tang, Xueming Wang, Jing Liu, Hanqing Lu. Multimedia News Summarization in Search. *ACM Trans. on Intelligent Systems and Technology* 7(3): 33 (1–20), 2016.
- Zechao Li, Jing Liu, Meng Wang, Changsheng Xu, Hanqing Lu. Enhancing News Organization for Convenient Retrieval and Browsing. *ACM Trans. on Multimedia Computing, Communications and Applications* 10(1): 1 (1–20), 2013.
- Zechao Li, Jing Liu, Xiaobin Zhu and Hanqing Lu: Multi-modal Multi-correlation Person-centric News Retrieval. In *ACM Conference on Information and Knowledge Management*, 2010.

Acknowledgements

First of all, I would like to express my gratitude to my supervisor, Prof. Hanqing Lu, for his long-term support and help. Professor Lu provided a well-equipped and active working environment for me and gave me full freedom to investigate any research problem of interest. His valuable suggestions and criticism play a significant role on my way toward a full-fledged researcher.

I would also thank Prof. Changsheng Xu (IEEE Fellow), for his guidance during my visiting study at China-Singapore Institute of Digital Media (CSIDM). Professor Jing Liu is another guider of my research. She gave me many pieces of advice in my research methodologies and her revision on my papers that I achieved remarkable progress. Her passion for academic research is of great importance to my decision to start my research career. I also thank Prof. Jian Cheng, Jinqiao Wang, Yifan Zhang, Meng Wang, Jinhui Tang, and Richang Hong. They spent a lot of time with me and provided me with a great deal of assistance. Much of my academic inspiration stems from discussions with them.

In addition, thanks are due to my colleagues in Image and Video Analysis (IVA) lab, including Chao Liang, Yang Liu (Male), Chunjie Zhang, Bo Wang, Chuanghua Gui, Xiao Yu Zhang, Si Liu, Tianzhu Zhang, Xiaobin Zhu, Peng Li, Yu Jiang, Yang Liu (Female), Wei Fu, Biao Niu, Yong Li, Cong Leng, Ting Yuan, among others, my colleagues in MultiMedia Computing (MMC) group, including Jitao Sang, Weiqing Min, Lei Yu, Zhaoquan Yuan, Ming Yan, among others. They helped me a lot during these years in many aspects. I will never forget the joyful and rewarding days spent with them. I am also grateful to all my good friends, wherever they are. My thanks go to everyone who contributed to my progress and happiness.

This book was partially supported by the 973 Program of China (Project No. 2014CB347600), and the National Natural Science Foundation of China under Grant No. 61402228.

Contents

1	Introduction	1
1.1	Multimedia Analysis and Understanding	1
1.2	Challenges and Progresses	3
1.3	Understanding-Oriented Multimedia Content Analysis	5
1.4	Organization of This Book	6
	References.	7
2	Understanding-Oriented Unsupervised Feature Selection	11
2.1	Introduction	11
2.2	Related Work	13
2.3	Clustering-Guided Sparse Structural Learning for Unsupervised Feature Selection.	14
2.3.1	Nonnegative Spectral Clustering	15
2.3.2	Sparse Structural Analysis.	16
2.3.3	Optimization	17
2.3.4	Convergence Analysis.	21
2.3.5	Discussions	22
2.3.6	Experiments	25
2.4	Nonnegative Spectral Analysis and Redundancy Control for Unsupervised Feature Selection.	32
2.4.1	The Objective Function.	32
2.4.2	Optimization	34
2.4.3	Experiments	37
2.5	Discussions.	40
	References.	44
3	Understanding-Oriented Feature Learning	47
3.1	Introduction	47
3.2	Related Work	48
3.2.1	Feature Selection.	48
3.2.2	Subspace Learning	49

- 3.3 The Proposed RSSL Framework. 50
 - 3.3.1 Formulation. 50
 - 3.3.2 Optimization 52
 - 3.3.3 Computational Complexity Analysis 56
 - 3.3.4 Image Understanding Tasks 56
- 3.4 Performance Evaluation 58
 - 3.4.1 Image Tagging 58
 - 3.4.2 Clustering 65
 - 3.4.3 Classification. 67
 - 3.4.4 Discussion. 69
- 3.5 Discussions. 71
- References. 71
- 4 Personalized Tag Recommendation 75**
 - 4.1 Introduction 75
 - 4.2 Related Work 77
 - 4.3 Personalized Geo-Specific Tag Recommendation 79
 - 4.3.1 Overview of Our Solution. 79
 - 4.3.2 Discovering Intermediate Space and Unified Space. 80
 - 4.3.3 Optimization 82
 - 4.3.4 Tag Representation in Unified Space 85
 - 4.3.5 Tag Recommendation 86
 - 4.4 Performance Evaluation 87
 - 4.4.1 Dataset 87
 - 4.4.2 Evaluation Measures 89
 - 4.4.3 Parameter Setting 90
 - 4.4.4 Compared Methods. 90
 - 4.4.5 Parameter Sensitivity Analysis 91
 - 4.4.6 Experimental Analysis. 94
 - 4.5 Discussions. 97
 - References. 97
- 5 Understanding-Oriented Multimedia News Retrieval 101**
 - 5.1 Introduction 101
 - 5.2 Related Work 103
 - 5.3 Person-Centered Multimedia News Retrieval 105
 - 5.3.1 Overview of the Proposed MMPNR System. 105
 - 5.3.2 Correlation Initialization 107
 - 5.3.3 Correlation Reconstruction 109
 - 5.3.4 Ranking and Visualization 114
 - 5.4 Multimedia News Contextualization 115
 - 5.4.1 Overview of GeoVisNews. 115
 - 5.4.2 News Document Geo-Localization 116
 - 5.4.3 Image Enrichment. 120
 - 5.4.4 Result Ranking and Visualization 124

- 5.5 Discussions. 126
- References. 127
- 6 Understanding-Oriented Multimedia News Summarization 131**
 - 6.1 Introduction 131
 - 6.2 Previous Work 133
 - 6.3 Multimedia News Summarization 135
 - 6.3.1 Hierarchical Topic Structure 135
 - 6.3.2 Multimedia Topic Representation 136
 - 6.3.3 User Interface Overview 142
 - 6.4 Performance Evaluation 143
 - 6.4.1 Data Sets and Experimental Settings 143
 - 6.4.2 Parameter Sensitiveness. 144
 - 6.4.3 Summarization Evaluations 146
 - 6.4.4 Interface Evaluations 149
 - 6.5 Discussions. 151
 - References. 151
- 7 Conclusion 155**
 - 7.1 Promising Topics 155
 - 7.2 The Prospects 156

Notations

Throughout this book, the lowercase italic letters (i.e., i, j, n , etc.) and the uppercase italic letters (i.e., A, B, M , etc.) denote scalars, while the bold uppercase characters (i.e., \mathbf{W}, \mathbf{X} , etc.) and the bold lowercase characters (i.e., \mathbf{a}, \mathbf{x} , etc.) are utilized to denote matrices and vectors, respectively. For any matrix \mathbf{A} , \mathbf{a}^i means the i -th column vector of \mathbf{A} , \mathbf{a}_i means the i -th row vector of \mathbf{A} , A_{ij} denotes the (i, j) -element of \mathbf{A} and $\text{Tr}[\mathbf{A}]$ is the trace of \mathbf{A} if \mathbf{A} is square. \mathbf{A}^T denotes the transposed matrix of \mathbf{A} . The Frobenius norm of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is defined as $\|\mathbf{A}\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n A_{ij}^2 = \text{Tr}[\mathbf{A}^T \mathbf{A}]$. The $\ell_{2,p}$ -norm ($p \in (0, 1]$) of \mathbf{A} is defined as

$$\|\mathbf{A}\|_{2,p} = \left(\sum_{i=1}^r \left(\sqrt{\sum_{j=1}^t A_{ij}^2} \right)^p \right)^{\frac{1}{p}} = \left(\sum_{i=1}^r \|\mathbf{a}_i\|_2^p \right)^{\frac{1}{p}}. \quad (1)$$

Note that in practice, $\|\mathbf{a}^i\|_2$ could be close to zero. For this case, we can follow the traditional regularization way and define $D_{ii} = \frac{1}{\|\mathbf{a}^i\|_2 + \varepsilon}$, where ε is very small constant. When $\varepsilon \rightarrow 0$, it is easy to verify that $\frac{1}{\|\mathbf{a}^i\|_2 + \varepsilon}$ approximates $\frac{1}{\|\mathbf{a}^i\|_2}$. Furthermore, let \mathbf{I}_m denote the identity matrix in $\mathbb{R}^{m \times m}$.

Chapter 1

Introduction

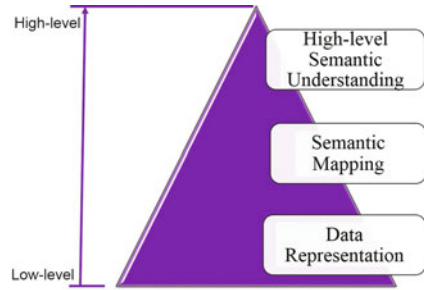
Abstract Multimedia content analysis is to understand the semantic information of multimedia data (such as text, image, audio, video, etc.). It is reasonable and necessary to develop understanding-oriented multimedia content analysis and incorporate model learning and understanding into a unified framework. In this chapter, we will first present an overview of multimedia content analysis and understanding, introduce the challenges and progresses in this field, and then describe the specifications of understanding-oriented multimedia content analysis. Finally, we give the organization of this book.

1.1 Multimedia Analysis and Understanding

With the popularity of intelligent devices (e.g., smart phones) and social media websites (e.g., flickr.com, youTube.com, etc.), multimedia data, especially images and videos, have been explosively increasing and playing an important role in our daily work and life. Taking facebook.com as an example, it reported in November 2013 that there are about 350 million images uploaded daily. There are 100h of video uploaded to YouTube every minute, resulting in an estimate of more than 2 billion videos totally by the end of 2013. This significantly extends the scope and application areas of multimedia. For example, Tencents free messaging and calling tool, Wechat, has attracted more than 300 million users in less than 2 years, which is tending to replace the traditional short message service (SMS). To say the least, we are really now living in a media world. Thereupon, it is necessary to develop approaches to intelligently analyze and understand the massive multimedia data.

Multimedia content analysis and understanding is to analyze and understand multimedia data (such as image, video, audio, graphic, etc.) using approaches from machine learning, artificial intelligence, and pattern recognition communities [7, 13]. In literatures, multimedia content analysis and understanding is deemed as a cross-disciplinary research area concerning with the intersection of image processing, computer vision, machine learning, artificial intelligence, pattern recognition, data mining, etc. It involves with techniques in visual/auditory physiology, signal processing, computer vision, information retrieval, etc. [43]. As shown in Fig. 1.1,

Fig. 1.1 Illustration of multimedia content analysis and understanding



there are usually three steps for multimedia content analysis and understanding, i.e., data representation, semantic mapping, and semantic applications. To better implement these three steps, it is necessary to analyze the characters of multimedia data.

- **Big** data. Multimedia is increasingly becoming the “biggest big data” as the most important and valuable source for insights and information. It covers from everyone’s experiences to everything happening in the world. As such, multimedia big data is spurring on tremendous amounts of research and development of related technologies and applications.
- **Social** data. Currently, multimedia data are mainly created and uploaded by users. Users can tag and comment on the uploaded multimedia data. Besides, there are rich metadata associated with multimedia data in the social sharing websites, such as EXIF, GPS, user, group, etc. That is, there are rich social context information associated with multimedia data.
- **Heterogeneous** data. Multimedia data contain multiple types of data, such as text, image, audio, video, etc. The same semantic can be described by data from different modalities. Essentially, data from different modalities have different properties, and they are heterogeneous. Besides, social multimedia data have more heterogeneous information, such as GPS information, user-provided tags, users’ comments, and the group information.

By considering the above characters, researchers have devoted substantial attention to developing methods for multimedia content analysis and understanding, such as data presentation [1, 18, 35, 40, 52, 55], metric learning [15, 21, 25, 39, 49], hashing [5, 30, 47, 48, 54], tag refinement and assignment [2, 22, 28, 50], semantic segmentation [26, 27, 32, 33, 36], and retrieval [6, 23, 45]. Thanks to the wide prevalence of multimedia data and the increasing demands for multimedia services, there has been a growing number of research on multimedia analysis and understanding, evidenced by both the volume of papers produced each year.

In spite of this, it is challenging to intelligently understand the multimedia content. It is well known that multimedia data are usually represented by the low-level features, such as images are described by the low-level visual features. There exists the well-known “Semantic Gap” [8], which is defined as the difference between the information that one can extract from the visual data and the interpretation that the

same data have for a user in a given situation. Besides, it is a challenge that how to uncover a better representation for multimedia data.

- **Multimedia Data Representation.** Multimedia data can be described by a variety of visual features, which are often quite different from each other. The dimension of data feature space is becoming increasingly large. It is inevitable to introduce noisy and/or redundant features. The effectiveness and efficiency of learning methods drop exponentially as the dimensionality increases, which is commonly referred to as the “curse of dimensionality” [16]. Therefore, it is a fundamental problem to find a suitable representation of high-dimensional data [4], which can enhance the performance of numerous tasks, such as multimedia analysis.
- **Semantic Mapping.** Due to the well-known semantic gap, it is challenging to identify better semantic mappings from the low-level feature space to the high-level semantic space. On the other hand, social multimedia data have rich context information, such as user-provided tags, users’ descriptions, users’ comments, GPS information, EXIF information, etc. It is beneficial for reducing the semantic gap to explore social context information, which can help to learn better semantic mappings.
- **Multimedia Understanding.** The ultimate goal is to understand the semantic information of multimedia data according to the corresponding application. With the proliferation of multimedia data, many interesting multimedia applications have been designed, such as mobile product retrieval, clothing retrieval [11, 29], multimedia news retrieval, recommendation [17, 37, 38], etc. Therefore, we should develop understanding-oriented approaches for the applications, which may help to attract more users and bring in bigger gains.

1.2 Challenges and Progresses

Although many previous methods have been proposed to address the problems of multimedia content analysis and understanding, there are some challenges for multimedia understanding and potential applications.

Multimedia Big Data. Multimedia is increasingly becoming the “biggest big dat,” which brings in new challenges. First, the traditional methods trained on a small-scale training set may cease to be effective for the multimedia big data. Second, for small data, it takes acceptable time to learn models. Unfortunately, it may take much time to deal with the multimedia big data. Consequently, how to efficiently and effectively deal with multimedia big data is important, and there is also an abysmal lack of new methods adaptively to multimedia big data.

Understanding-oriented Representation. Traditional features of multimedia data are extracted independent from the follow-up understanding tasks. From the perspective of representation, the features of multimedia data are always noisy and/or redundant, and the dimension of features is becoming increasingly high. For the hand-crafted features, people manually design a feature extraction pipeline by the