

ADVANCES ON
INFORMATION
PROCESSING AND
MANAGEMENT

CONFERENCE ORGANIZERS INSTITUTES

The International Conference on Integrated Information is supported by the following Institutes:

Emerald Group Publishing Limited
Technological educational Institute of Athens, Greece
University of Peloponnese, Greece
National And Kapodistrian University of Athens, Greece
Mednet Hellas, The Greek Medical Network
2nd AMICUS Workshop

To learn more about I-DAS, including the Book Series, please visit the webpage
<http://www.i-das.org/>

INTEGRATED INFORMATION

International Conference on Integrated Information

Kos, Greece September, 29 – October, 3 2011

EDITORS

Georgios A. Giannakopoulos
Technological Educational Institute of Athens, Greece

Damianos P. Sakas
University of Peloponnese, Greece

All papers have been peer-reviewed



Piraeus, Greece, 2011

Editors

Georgios A. Giannakopoulos

Technological Educational Institute of Athens
Faculty of Management and Economics
Department of Library Science and Information Systems
Address: Aghiou Spyridonos Street, 12210, Egaleo
E-mail: gian@teiath.gr

Damianos P. Sakas

University of Peloponnese
Faculty of Science and Technology
Department of Computer Science and Technology
Address: End of Karaiskaki St., 22100, Tripolis, Greece
E-mail: D.Sakas@uop.gr

The copyrights will be owned by the authors under the Creative Commons Attribution-Non Commercial license (<http://creativecommons.org/licenses/by-nc/3.0/>), which permits unrestricted use, distribution, and reproduction in any non commercial medium, provided the original work is properly cited.

ISSN:

Printed in the Greece, EU

CONTENTS

PREFACE: Proceedings of the International Conference on Integrated Information (IC-INFO 2011)	1
Georgios A. Giannakopoulos, Damianos P. Sakas	
Conference Details	3
Keynote Speaker	5
SYMPOSIUM ON INFORMATION AND KNOWLEDGE MANAGEMENT	6
Prof. Christos Skourlas	
Towards the Preservation and Availability of Historical Books and Manuscripts: A Case Study	8
Eleni Galiotou	
An Extensive Experimental Study on the Cluster-based Reference set Reduction for Speeding-up the k-nn Classifier	12
Stefanos Ougiaroglou, Georgios Evangelidis and Dimitris A. Dervos	
Exploiting the Search Culture Modulated by the Documentation Retrieval Applications	16
Nikitas N. Karanikolas and Christos Skourlas	
Information and Knowledge Organization: The Case of the TEI of Athens	22
Anastasios Tsolakidis, Manolis Chalaris and Ioannis Chalaris	
Providing Access to Students with Disabilities and Learning Difficulties in Higher Education through a Secure Wireless framework	26
Catherine Marinagi and Christos Skourlas	
Improving Query Efficiency in High Dimensional Point Indexes	30
Evangelos Outsios and Georgios Evangelidis	
Text Segmentation Using Named Entity Recognition and co-Reference Resolution in Greek Texts	34
Pavlina Fragkou	
KINISIS, a Graphical XQuery Language	42
Euclid Keramopoulos, Achilleas Pliakas, Konstantinos Tsekos and Ignatios Deligiannis	
Dimensionality Curse, Concentration Phenomenon and the KDB-tree	46
Nikolaos Kouiroukidis and Georgios Evangelidi	

Applying Balanced Scorecard Strategic Management in Higher Education	50
Manolis Chalaris, Anastasios Tsolakidis and Ioannis Chalaris	
A Web Portal Model for NGOs' Knowledge Management	54
Zuhal Tanrikulu	
The Digital Archives System and Application Optimized for the Tradition Knowledge Archives	58
Jeon Hong. Chan, In Deok. Hwang, Jae Hak. Park, Hyeok. Sim, U won. Gwon and Soon Cheol. Park	
A Semi-automatic Emerging Technology Trend Classifier Using SCOPUS and PATSTAT	62
Seonho Kim, Woondong Yeo, Byong-Youl Coh, Waqas Rasheed, Jaewoo Kang	
Presenting a Framework for Knowledge Management within a Web Enabled Living Lab	66
Lizette de Jager and Albertus AK Buitendag and Potjie (JS) van der Walt	
4TH SYMPOSIUM ON BUSINESS AND MANAGEMENT AND DYNAMIC SIMULATION MODELS SUPPORTING MANAGEMENT STRATEGIES	71
Dr. Damianos Sakas	
New Political Communication Practices: No Budget Events Management. The New Challenge	73
Evangelia N. Markaki, Damianios P. Sakas and Theodore Chadjipantelis	
Free Software – Open Source Software. A Powerful Tool for Developing Creativity in the Hands of the Student	78
Nasiopoulos K. Dimitrios, Damianos P. Sakas, Konstantinos Masselos	
Open Source Web Applications. How it Spread Through the Internet and their Contribution to Education.	82
Nasiopoulos K. Dimitrios, Damianos P. Sakas, Konstantinos Masselos	
Culture in Modern Times in the Frame of Luhmann's System Theory	85
Anastasia J. Chournazidis	
Managing Scientific Journals: A Cultural Viewpoint	87
Marina C. Terzi, Damianos P. Sakas, and Ioannis Seimenis	
A Conceptual Framework for Analyzing Knowledge-based Entrepreneurship	92
Nikos S. Kanellos	

SESSION ON INFORMATION HISTORY: PERSPECTIVES, METHODS AND CURRENT TOPICS	96
Prof. Laszlo Karvalics	
Emerging Research Fields in Information History	98
Laszlo Z. Karvalics	
Information Management through Elementary Data Clusters: New Observations on Pridianum-Type Roman Statistical Documents	102
Gergő Gellérfi	
Information and Secrecy on the Silk Road. Methods of Encryption of Legal Documents in Inner Asia (3th-4th century)	106
Szabolcs Felföldi	
The Role of Information and Disinformation in the Establishment of the Mongolian Empire: A Re-examination of the 13th century Mongolian History from the Viewpoint of Information History	110
Márton Gergő Vér	
Early Warning Systems and the Hospitallers in the Eastern Mediterranean	114
Zsolt Hunyadi	
Information Management as Establishment Dutch Navigational Knowledge on Japan, 1608-1641	118
Gabor Szommer	
Files Everywhere - Register and Training of Men for Military and Civil Purpose in Prussia in the early 18th century	123
Marton Holczer	
SYMPOSIUM ON INTEGRATED INFORMATION: THEORY, POLICIES, TOOLS	126
Prof. Georgios Giannakopoulos	
Approaching Information as an Integrated Field: Educating Information Professionals	128
Georgios Giannakopoulos, Daphne Kyriaki Manesi and Stryidon Zervos	
Special Libraries as Knowledge Management Centers	132
Eva Semertzaki	
Digital Libraries' Developers and their Suitability: A Case Study	136
Maria Monopoli	

A Preliminary Study for the Creation of a Greek Citation index in the Humanities and the Social Sciences (GCI – H&SS)	140
Daphne Kyriaki-Manessi and Evi Sachini	
Archiving as an Information Science. Evidence from a Survey Carried out on a Sample of Greek Students	144
Georgios Giannakopoulos and Ioannis Koumantakis	
Transition Process of E-records Management and Archiving System in Universities: Ankara University	147
S. Özlem Bayram and Fahrettin Ozdemirci	
Government Information: Access and Greece's Efforts for Access	150
Aikaterini Yiannoukakou	
School Archives and their Potentials in Teaching: Aspects of Greek Reality	156
Sonia Geladaki and Panagiota Papadimitriou	
Research on School Libraries in Greece and Suggestions on its Further Development	160
Georgios D. Bikos	
Building Digital Collections for Archeological Sites: Metadata Requirements and CIDOC CRM Extension	164
Georgios S. Gkrous and Mara Nikolaidou	
Museological Claims to Autonomous Knowledge: Rethinking the Conceptual Mode of Display and its Claims to Knowledge	169
Assimina Kaniari and Georgios Giannakopoulos	
Use of Library Loan Records for Book Recommendation	172
Keita Tsuji, Erika Kuroo, Sho Sato, Ui Ikeuchi, Atsushi Ikeuchi, Fuyuki Yoshikane and Hiroshi Itsumura	
Developing a National Database on Librarianship and Information Science. The Case of E-VIVA, the Hellenic Fulltext Database	176
Filippos Ch. Tsimpoglou, Vasiliki V. Koukounidou and Eleni K. Sakka	
Integrated Access to Cultural Heritage Information Pieces in Iran Astan-Quds Razavi's Organization of Libraries, Museums and Documents Center: A Theory of Unionization Disparate Information Assets over Imam Reza's Zarih	181
Ms. Mitra Zarei and Ms. Maliheh Farrokhnia	
Attitudes of University Librarians and Information Scientists towards the Draft Code of	185

Library Ethics to Present a Model for Final Library Ethical Codes	
Mahsoomeh Latifi, Fatemeh Zandian and Hasan Siamian	
SESSION ON OPEN ACCESS REPOSITORIES: SELF-ARCHIVING, METADATA, CONTENT POLICIES, USAGE	188
Dr. Alexandros Koulouris	
Geographical Collections in Greek Academic Libraries: Current Situation and Perspectives	189
Ifigenia Vardakosta and Sarantos Kapidakis	
Information Seeking Behavior: Factors that Affect the Behavior of Greek Astronomers	194
Hara Brindesi and Sarantos Kapidakis	
Aggregating Metadata for Europeana: The Greek Paradigm	198
Alexandros Koulouris, Vangelis Banos and Emmanouel Garoufallou	
Integrating a Repository with Research Output and Publications: The Case of the National Technical University of Athens	202
Dionysis Kokkinos	
Implementation of Workflows as Finite State Machines in a National Doctoral Dissertations Archive	205
Nikos Houssos, Dimitris Zavaliadis, Kostas Stamatis and Panagiotis Stathopoulos	
Practices of “Local” Repositories of Legally Protected Immovable Monuments. A Global Scheme for ‘Designation – Significance’ Information	209
Michail Agathos and Sarantos Kapidakis	
Integration of Metadata in BWMETA-2.0.0 Format	213
Katarzyna Zamlynska, Jakub Jurkiewicz and Lukasz Bolikowski	
SESSION ON EVIDENCE-BASED INFORMATION IN CLINICAL PRACTICE	216
Dr. Evangelia Lappa	
Applicability of Data Mining Algorithms on Clinical Datasets	218
Wilfred, Bonney	
Changing Roles of Health Librarians with Open Access Repositories	221
Christine Urquhar and Assimina Vlachaki	
From Medical Records to Health Knowledge Management Systems: The Coding to Health Sector	225
Evangelia C. Lappa and Georgios A. Giannakopoulos	

The Survey of Skill, Attitude and Use of Computer and Internet among Faculty Members	229
Hasan Siamian, Azita Bala Ghafari, Kobra Aligolbandi, Mohammad Vahedi and Gholam Ali Golafshani Jooybari	
Trends in Scholarly Communication among Biomedical Scientists in Greece	232
Assimina Vlachaki and Christine Urquhart	
SESSION ON ELECTRONIC PUBLISHING: A DEVELOPING LANDSCAPE	236
Dr. Dimitris Kouis	
E-Journal and Open Access Journal Publishing in the Humanities: Preliminary Results from a Survey among Byzantine Studies Scholars	238
Victoria Tsoukala and Evi Sachini	
Preliminary Results on a Printed VS Electronic Text Books Assessment Through Questionnaire	242
Dimitrios A. Kouis and Kanella Pouli	
An Interpretation of Aristotelian Logic According to George Boole	246
Markos N. Dendrinis	
SESSION ON INFORMATION CONTENT PRESERVATION AS OUTCOME OF CONSERVATION OF CULTURAL HERITAGE: ETHICS, METHODOLOGY AND TOOLS	251
Prof. George Panagiaris and Dr. Spiros Zervos	
Intrinsic Data Obfuscation as the Result of Book and Paper Conservation Interventions	254
Spiros Zervos, Alexandros Koulouris and Georgios Giannakopoulos	
Mass Deacidification: Preserving More than Written Information	258
Michael Ramin, Evelyn Eisenhauer and Markus Reist	
Information Literacy of Library Users: A Case Study of Mazandaran Public Library Users, Iran	260
Hussein Mahdizadeh and Hasan Siamian	
The Narratives of Paper in The Archives of the New Independent Greek State (Mid 19th c.)	264
Ourania Kanakari and Maria Giannikou	
From Macro to Micro and from Micro to Nano: The Evolution of the Information Content Preservation of Biological Wet Specimen Collections	268
Nikolaos Maniatis and Georgios Panagiaris	

Digital images: A valuable scholar's tool or misleading material?	272
Patricia Engel	
Attitudes of University Librarians and Information Scientists Towards the Draft Code of Library Ethics to Present a Model for Final Library Ethical Codes	277
Mahsoomeh Latifi, Fatemeh Zandianand and Hasan Siamian	
Investigation of the Degradation Mechanisms of Organic Materials: From Accelerated Ageing to Chemometric Studies	280
Ekaterini Malea, Effie Papageorgiou and Georgios Panagiaris	
SESSION ON DIVERGENCE AND CONVERGENCE: INFORMATION WORK IN DIGITAL CULTURAL MEMORY INSTITUTIONS	285
Dr. Susan Myburgh	
Extending Convergence and Divergence in Cultural Memory Institutions: The Old Slave Lodge in the New South Africa	287
Archie L Dick	
The Transfer of Knowledge from Large Organizations to Small: Experiences from a Research Project on Digitization in Wales	289
Clare Wood-Fisher, Richard Gough, Sarah Higgins, Menna Morgan, Amy Staniforth and Lucy Tedd	
The Usage of Reference Management Software (Rms) in an Academic Environment : A Survey at Tallinn University	293
Enrico Francese	
Varialog : How to Locate Words in a French Renaissance Virtual Library	297
Marie-Hélène Lay	
The Urge to Merge: A Theoretical Approach	301
Susan Myburgh	
SYMPOSIUM ON ADVANCES INFORMATION FOR STRATEGIC MANAGEMENT	304
Professor Nikolaos Konstantopoulos	
Empowerment in the Tax Office of Greece	306
Antonios E. Giokas and Nikolaos P. Antonakas	
Building Absorptive Capacity Through Internal Corporate Venturing	310
Ioannis M. Sotiriou and Alexandros I. Alexandrakis	

The Monitoring Information System (M.I.S.) - An information and Management System for Projects Co-financed Under the National Strategic Reference Framework (NSRF) and the Community support framework (CSF)	314
Catherina G. Siampou, Eleni G. Fassou and Athanassios P. Panagiotopoulos	
Corruption in Tax Administration: The Entrepreneurs View Point	318
Nikolaos P. Antonakas, Antonios E. Giokas and Nikolaos Konstantopoulos	
Conflicts between the IT Manager and the Software House after the Strategic Choice of Outsourcing of the Information Processes in Maritime Companies.	322
Anthi Z. Vaxevanou, Nikolaos Konstantopoulos, Damianos P. Sakas	
Contemporary Forms of Ordering Between the Supply Department and Ship Chandler Companies in the Shipping Industry	325
Anthi Z. Vaxevanou, Nikolaos Konstantopoulos, Damianos P. Sakas	
Strategies Implemented and Sources Used for the Acquisition of Information on Foreign Markets	329
Myropi Garri, Nikolaos Konstantopoulos and Michail G. Bekiaris	
The Effect of High Performance Working Systems on Informative Technology in Enterprises after Organisation Changes such as Mergers & Acquisitions	333
Nikolaos Konstantopoulos and Yiannis Triantafyllopoulos	
Personnel's Absorptive Capacity as a Guiding Concept for Effective Performance in Informative Technology	337
Nikolaos Konstantopoulos and Yiannis Triantafyllopoulos	
SESSION ON CONTEMPORARY ISSUES IN MANAGEMENT: ORGANISATIONAL BEHAVIOUR, INFORMATION TECHNOLOG, EDUCATION & HOSPITAL LEADERSHIP	341
Dr. Panagiotis Trivellas	
Investigating the Importance of Sustainable Development for Hotel SMES	343
Panagiotis Reklitis and Anestis Fotiadis	
Strategic Alignment of ERP, CRM and E-business: A Value Creation	347
Catherine C. Marinagi and Christos K. Akrivos	
The Impact of Occupational Stress on Performance in Health Care	351
Panagiotis Trivellas Panagiotis Reklitis and Charalambos Platis	

The Impact of Emotional Intelligence on Job Outcomes and Turnover Intention in Health Care	356
Panagiotis Trivellas Vassilis Gerogiannis and Sofia Svarna	
SYMPOSIUM ON BUSINESS MANAGEMENT AND COMMUNICATION STRATEGIES SUPPORTING DECISION MAKING PROCESS IN TOURISM SECTOR	360
Dr. Panagiota Dionysopoulou	
The Human Factor as a Mediator to the Total Quality in the Tourism Companies. The impact of Employees' Motivation to Quality Improvements	362
Christos K. Akrivos and Panagiotis Reklitis	
Tourist Destination Marketing and Management Using Advanced ICTS Technologies	365
Anastasia Argyropoulou, Panagiota Dionyssopoulou, Georgios Miaoulis	
G.N.T.O. (Greek National Tourism Organization) Communication Strategy in Advertising Campaigns 1991-2006	370
George Stafylakis and Panagiota Dionyssopoulou	
GENERAL PAPERS	375
The role of Environmental Education within the Framework of the Environmental Policy of a Regional Municipality	376
Vassiliki Delitheou and Dimitra Thanasia	
Issues of Social Cohesion: A case study from the Greek Urban Scenery	380
Evgenia Tousi	
Merging Activity and Employee Performance: The Greek Banking System	384
Panagiotis Liargovas and Spyridon Repousis	
Sustainable Development and Corporate Social Responsibility in Higher Education: Some Evidence from Greece	387
Anastasios Sepetis and Fotios Rizos	
Exploring the Effects of Organizational Culture on Collaborative vs. Competitive Knowledge Sharing Behaviors	395
Hanan Abdulla Mohammed Al Mehairi and Norhayati Zakaria	

Preface: Proceedings of the International Conference on Integrated Information (IC-ININFO 2011)

GEORGIOS A. GIANNAKOPOULOS

Department of Library Science and Information Systems, Technological Educational Institute of Athens, Aghiou Spyridonos, Egaleo, 12210, Greece

DAMIANOS P. SAKAS

Department of Computer and Technology Science, University of Peloponnese, Praxitelous 89-91, Piraeus, 18532, Greece

Aims and Scope of the Conference

The International Conference on Integrated Information 2011 took place in Kos Island, Greece, between September, 29 and October, 3, 2011. IC-ININFO is an international interdisciplinary conference covering research and development in the field of information management and integration.

The conference aims at creating a forum for further discussion for an Integrated Information Field incorporating a series of issues and/or related organizations that manage information in their everyday operations. Therefore, the call for papers is addressed to scholars and/ or professionals of the fields of Library and Archives Science (including digital libraries and electronic archives), Museum and Gallery Studies, Information Science, Documentation, Information Management, Records Management, Knowledge Management, Data management and Copyright experts the latter with an emphasis on Electronic Publications. Furthermore, papers focusing on issues of Cultural Heritage Management and Conservation Management are also be welcomed along with papers regarding the Management of Nonprofit Organizations such as libraries, archives and museums.

One of the primary objectives of the IC-ININFO will be the investigation of information-based managerial change in organizations. Driven by the fast-paced advances in the Information field, this change is characterized in terms of its impact on organizations that manage information in their everyday operations.

Grouping emerging technologies in the Information field together in a close examination of practices, problems and trends, IC-ININFO and its emphases on integration and management will present the state of the art in the field. Addressed jointly to the academic and practitioner, it will provide a forum for a number of perspectives based on either theoretical analyses or empirical case studies that will foster dialogue and exchange of ideas.

Topics of general Interest

Library Science, Archives Science, Museum and Gallery Studies, Information Science, Documentation, Digital Libraries, Electronic Archives, Information Management, Records / Document Management, Knowledge Management, Data Management, Copyright, Electronic Publications, Cultural Heritage Management, Conservation Management, Management of Nonprofit Organizations, History of Information, History of Collections, Health Information

Symposia

The Conference offered a number of sessions under its patronage, providing a concise overview of the most current issues and hands-on experience in information-related fields.

- Symposium on Integrated information: Theory, Policies, Tools
- 4th Symposium on Business and Management and Dynamic Simulation Models supporting management strategies

- Session on Open Access Repositories: Self-archiving, Metadata, Content policies, Usage
- Session on Evidence-Based Information in Clinical Practice
- Session on Business Management and Communication Strategies supporting Decision Making Process in Tourism Sector
- Session on Electronic Publishing: A Developing Landscape
- Session on Information and Knowledge Management
- Session on Information Content Preservation as Outcome of Conservation of Cultural Heritage: Ethics, Methodology and Tools
- Session on Advances Information for Strategic Management
- Session on Information History: Perspectives, Methods and Current Topics
- Session on Divergence and Convergence: Information Work in Digital Cultural Memory Institutions
- Session on Contemporary issues in Management: Organisational Behaviour, Information Technology, Education & Hospital leadership.

The wide range of aspects that the sessions covered, highlighted future trends in the Information Science.

Paper Peer Review

More than 300 papers had been submitted for consideration in IC-ININFO 2011. From them, 91 were selected for presentation, after peer review in a double blind review process. The accepted papers were presented at IC-ININFO 2011.

Thanks

We would like to thank all members that participated in any way in the IC-ININFO 2011 Conference and especially:

- The famous publishing house Emerald for its communication sponsorship.
- The co-organizing Universities and Institutes for their support and development of a high-quality Conference scientific level and profile.
- The members of the Scientific Committee that honored the Conference with their presence and provided a significant contribution to the review of papers as well as for their indications for the improvement of the Conference.
- All members of the Organizing Committee for their help, support and spirit participation before, during and after the Conference.
- The Session Organizers for their willing to organize sessions of high importance and for their editorial work, contributing in the development of valued services to the Conference.
- PhDC Marina Terzi for her excellent editorial work, contributing in the production of the Conference proceedings.

CONFERENCE DETAILS

Chairs

Georgios A. Giannakopoulos, Technological Educational Institute of Athens, Greece
Damianos P. Sakas, University of Peloponnese, Greece

Co-Chairs

Daphne Kyriaki – Manesi, Technological Educational Institute of Athens, Greece
Dimitrios Vlachos, University of Peloponnese, Greece

Scientific Committee

Amanda Spink, Queensland University of Technology
Andreas Bagias, European Court
Andreas Rauber, Vienna University of Technology
Astrid van Wesenbeeck, SPARC Europe
Christine Urquhart, Aberystwyth University
Christos Schizas, University of Cyprus
Christos Skourlas, Technological Educational Institute of Athens
Claire Farago, University of Colorado at Boulder
Claus-Peter Klas, FernUniversität in Hagen
Costas Vassilakis, University of Peloponnese,
Dimitris Dervos, Technological Educational Institute of Thessaloniki
Eelco Ferwerda, OAPEN
Elena Garcia Barriocanal, University of Alcalá
Emmanuel Garoufallou, Technological Educational Institute of Thessaloniki
Filippos Tsimpoglou, University of Cyprus
Fillia Makedon, University of Texas at Arlington
George Korres, University of Newcastle
Georgios Evangelidis, University of Macedonia
Georgios Panagiaris, Technological Educational Institute of Athens
Johan Oomen, Netherlands Institute for Sound and Vision
José Aldana, University of Malaga
Konstantinos Masselos, University of Peloponnese
Luciana Duranti, The University of British Columbia
Markos N. Dendrinis, Technological Institute of Athens
Milena Dobрева, University of Strathclyde
Prodromos Tsiavos, London School of Economics and Political Science
Sándor Darányi, University of Borås
Sarantos Kapidakis, Ionian University
Sirje Virkus, Tallinn University
Spiros Zervos, Technological Educational Institute of Athens
Susan Myburgh, University of South Australia
Theodoros Pitsios, University of Athens, Faculty of Medicine

Organizing Committee

Alexandros Koulouris (Chair), Technological Educational Institute of Athens
Christos Christopoulos, SCEV Scientific Events Ltd
Marina Terzi, University of the Aegean, Greece
Evangelia Markaki, Aristotle University of Thessaloniki

Assimina Kaniari, Athens School of Fine Arts
Evangelia Lappa, General Hospital Attikis K.A.T.
Dimitris Kouis, Greek Ministry of Education, Lifelong Learning and Religious Affairs
Dionysis Kokkinos, National Technical University of Athens

KEYNOTE SPEAKER



Professor Amanda Spink

Professor Amanda Spink has published over 340 scholarly journal articles, refereed conference papers and book chapters, and 6 books. Many of her journal articles are published in the *Journal of the American Society for Information Science and Technology*, *Information Processing and Management*, and the *Journal of Documentation*. She is Editor of the Emerald journal *Aslib Proceedings*. Amanda's research has been published at many conferences including ASIST, IEEE ITCC, CAIS, Internet Computing, ACM SIGIR, and ISIC Conferences. Her recent books include *Information Behavior: An Evolutionary Instinct* and *Web Search: Multidisciplinary Perspectives*, both published by Springer. Amanda's research focuses on theoretical and empirical studies of information behavior, including the evolutionary and developmental foundations. The National Science Foundation, the American Library Association, Andrew R. Mellon Foundation, Amazon.com, Vivisimo. Com, Infospace.com, NEC, IBM, Excite.com, AlltheWeb.com, AltaVista.com, FAST, and Lockheed Martin have sponsored her research. In 2008 Professor Spink had the second highest H-index citation score in her field from 1998 to 2008 [Norris, M. (2008)]. Ranking Fellow Scholars and their H-Index: Preliminary Survey Results. Loughborough University, Dept of Information Science Report].

Dimensionality Curse, Concentration Phenomenon and the KDB-tree

Nikolaos Kouiroukidis and Georgios Evangelidis

University of Macedonia, Department of Applied Informatics, 54006, Thessaloniki, Greece
 {kouiruki, gevan} (at) uom.gr

Abstract: *The problem of indexing large volumes of high dimensional data is an important and popular issue in the area of database management. There are many indexing methods that behave well in low dimensional spaces, but, in high dimensionalities, the phenomenon of the curse of dimensionality renders all indexes useless. For example, when issuing range queries almost all of the index pages have to be retrieved for answering these queries. In this paper we review the state-of-the-art research regarding high dimensional spaces and we demonstrate the dimensionality curse phenomenon using the TPIE KDB-tree implementation.*

Keywords: *Dimensionality curse, KDB tree, Hypercube range queries*

I. INTRODUCTION

The term “curse of dimensionality” describes the rapid deterioration in the performance of high dimensional indexes as the number of variables (or dimensions) increases. When range or k-nearest neighbor queries are issued in high dimensional spaces, most (if not all) of the pages of the indexing structures that are employed to store the high dimensional points are visited, and the good performing in low dimensional spaces indexing methods, end up behaving as the plain sequential scan.

One of the classical indexing methods is the KDB-tree (Robinson, 1981) with TPIE (Arge et al, 2002) being one of his most efficient implementations. The KDB-tree combines some of the properties of the adaptive k-d-tree (Bentley, 1975) and the B-tree to handle multidimensional points. Each interior node corresponds to an interval-shaped region. Regions corresponding to nodes at the same tree level are mutually disjoint; their union is the complete universe. The leaf nodes store the data points that are located in the corresponding partition. Like the B-tree, the KDB-tree is a perfectly balanced tree that adapts well to the distribution of data.

In Section II, we present some observations regarding the dimensionality curse phenomenon. In Section III, we discuss the concentration phenomenon and in Section IV, we demonstrate the behavior of the KDB-tree in high dimensions. We conclude in Section V.

II. THE CURSE OF DIMENSIONALITY

The following phenomena give an insight to the notion of the dimensionality curse. See Weber et al. (1998) for further details.

1. The partitioning schemes usually split the data space in each dimension in two halves. With d dimensions there are 2^d partitions. With $d \leq 10$ and N on the order of 10^6 such a partition makes sense. However if d is larger, say $d=100$, there are around 10^{30} partitions for only 10^6 points. An overwhelming number of partitions are empty.

2. If we consider a hypercube range query with length s in all d dimensions the probability that a point lies within that range query is given by $P^d[s]=s^d$. This probability function is plotted in Fig. 1 below. From the formula, directly follows that even very large range queries are not likely to contain a point. At $d=100$ a range query with length 0.95 selects 0,59% of the data points. This hypercube range query can be placed anywhere in the data space Ω . Thus, we conclude that the data space is sparsely populated.

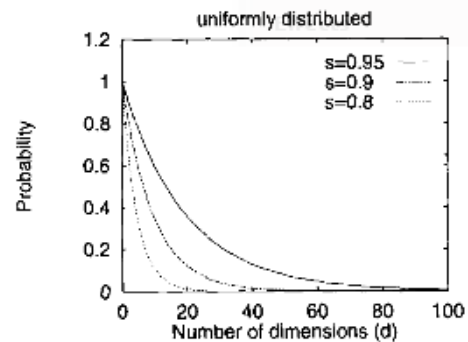


Figure 1. Plotting the probability that a hypercube query with side s contains a point.

3. The largest spherical query that fits entirely within the data space is the query $sp^d(Q, 0.5)$ where Q is the centroid of the data space. The probability that an arbitrary point R lies within this sphere is given by the sphere volume

$$P\{R \in sp^d(Q, \frac{1}{2})\} = \frac{Vol(sp^d(Q, \frac{1}{2}))}{Vol(\Omega)} = \frac{\sqrt{\pi^d} \cdot (\frac{1}{2})^d}{\Gamma(\frac{d}{2} + 1)}$$

The relative volume of the sphere shrinks markedly as the dimensionality grows and it increasingly becomes improbable that any point will be found within this sphere at all. Table 1 shows this probability for various numbers of dimensions.

4. From the probability equation given above, one can determine a size a data set would have to have

such that on average at least one point falls into the sphere $sp^d(Q,0.5)$ (for even d). This is given in the following equation:

$$N(d) = \frac{(\frac{d}{2})!}{\sqrt{\pi^{\frac{d}{2}}} \cdot (\frac{1}{2})^{\frac{d}{2}}}$$

Table 1 enumerates this function for various numbers of dimensions. The number of points needed explodes exponentially. At $d=20$, a database must contain at least 40 million points in order to ensure that on average at least one point lies within this sphere.

D	P[R \in $sp^d(Q,0.5)$]	N(d)
2	0.785	1.273
4	0.308	3.242
10	0.002	401.5
20	$2.461 * 10^{-8}$	40,631,627
40	$3.278 * 10^{-21}$	$3.050 * 10^{20}$
100	$1.868 * 10^{-70}$	$5.353 * 10^{69}$

Table 1. Probability that a point is in the largest hyper-sphere

5. The expected Nearest Neighbor distance between two points in a data space Ω is given by the following formula

$$E[nn^{dist}] = \int_{Q \in \Omega} E[Q, nn^{dist}] dQ$$

where Q is the query point. Based on this formula, and if one estimates it with the Monte Carlo method, one finds that NN distance grows steadily with d , and except trivially small data sets, the objects are widely scattered and the probability of being able to identify a good partitioning of the data space diminishes.

6. Finally, due to the dimensionality curse phenomenon, as we will demonstrate in our experiments with the KDB-tree, when a range query is performed nearly all data pages have to be accessed in order to obtain the answer. This equals almost to a sequential scan.

III. CONCENTRATION PHENOMENON

The concentration phenomenon can be stated as follows (Ledoux, 2001): in high dimensional spaces all pairwise distances between points seem identical. Here, we'll study the concentration of the distances through the concentration of the norm. If we have n points with d dimensions each, taking values from the unit cube $[0,1]^d$ and we then consider their norms $\|x\|$, the values of $\|x\|$ are bounded in the interval $[0,M]$, where $M=\|(1,1,\dots,1)\|$.

Let us consider the euclidean norm $M=\sqrt{d}$. If we plot the minimum observed value and the maximum observed value, we observe that in low

dimensions these values are close to the bounds of the domain of the norm, respectively 0 and \sqrt{d} . Also, the average value of the norm increases with the dimension, whereas the standard deviation seems rather constant. When the dimension is large (above 10) the minimum and maximum observed values tend to move away from the bounds. When the number of points are, for example, 100000 all the observed norms seem to concentrate in a small portion of their domain. In addition this portion gets smaller and smaller as the dimension grows when compared to the size of the total domain.

The Minkowski norms form a family of norms parametrized by their exponent $p=1,2,3,\dots$

$$\|X\|_p = \left(\sum_i |X_i|^p \right)^{\frac{1}{p}}$$

When $0 < p < 1$, the triangle inequality does not hold so these norms are called prenorms or fractional norms. Actually, the inequality is reversed. A consequence is that the straight line is no longer the smallest path between two points. Fig. 2 depicts 2D unit balls (that is the set of x^j for which $\|x^j\|=1$) for various values of p . We see that for $p \geq 1$ the balls are convex and for $0 < p < 1$ they are not.

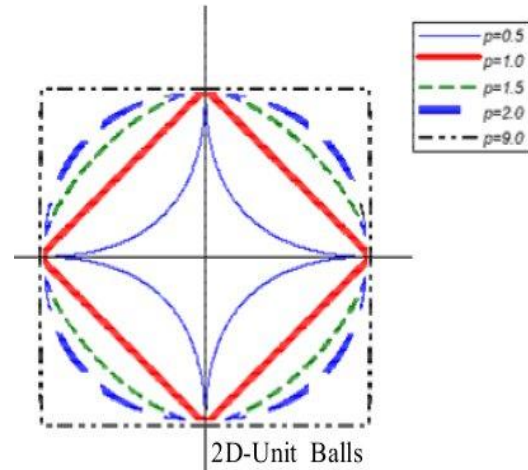


Figure 2. 2D-Unit Balls.

A. Concentration of the euclidean norm

If X is in R^d and is a random vector with independent and identically distributed components, and X_i follows distribution F , then

$$E(\|X\|_2) = \sqrt{ad - b} + O(1/d) \text{ and}$$

$$\text{Var}(\|X\|_2) = b + O(1/\sqrt{d}), \quad \text{where } a \text{ and } b$$

are constants that do not depend on the dimension (François et al., 2007; Aggarwal, 2001). This holds for

any kind of distribution. Different distributions will lead to different values for a and b but the asymptotic results remain.

This shows that the expectation of the euclidean norm of random vectors increases as the square root of the dimension, whereas its variance is constant and independent of the dimension. Therefore, when the dimension is large the variance of the norm is very small compared with its expected value. Also when the dimension is large vectors seem normalized. The relative error made while considering $E(\|X\|_2)$ instead of the real value of $\|X\|_2$ becomes negligible. As a consequence, high dimensional vectors appear to be distributed on a sphere of radius $E(\|X\|_2)$.

Since the euclidean distance is the norm of the difference between two random vectors, it's expectation and variance follow the two above laws and pairwise distances between points in high dimensional spaces seem to be all identical. Finally, if X_i are not independent the results are still valid provided that we replace d with the actual number of degrees of freedom.

In contrast to the work of Demartines (1994), where a data set X consists of n independent draws x^j from a single random vector X, Beyer (1999) considers n random vectors P^j where a dataset is made of one realization of each random vector. Beyer's theorem states that if P^j $1 \leq j \leq n$ are n d-dimensional independent and identically distributed random vectors and if

$$\lim_{d \rightarrow \infty} \text{Var} \left(\frac{\|P^{(j)}\|}{E(\|P^{(j)}\|)} \right) = 0$$

then for any $\epsilon > 0$

$$\lim_{d \rightarrow \infty} \mathbf{P} \left[\frac{\max_j \|P^{(j)}\| - \min_j \|P^{(j)}\|}{\min_j \|P^{(j)}\|} \leq \epsilon \right] = 1.$$

This is explained as follows. Suppose there are a set of n data points randomly distributed in the d-dimensional space and some query points are supposed to be located at the origin without loss of generality. Then, if the above hypothesis is satisfied, independent of the distribution of the components of the P_j , the difference between the largest and smallest distances to the query point becomes smaller and smaller when compared with the smallest distance when the dimension increases. The ratio

$$\frac{\max_j \|P^{(j)}\| - \min_j \|P^{(j)}\|}{\min_j \|P^{(j)}\|}$$

is called the relative contrast.

So, Beyer concluded that all points are located at approximately the same distance from the query

point. Thus, the concept of NN in a high dimensional space is less intuitive than in a lower dimensional one.

B. Concentration of Minkowski norms

There is the theorem of Hinneburg (François et al., 2007; Aggarwal et al., 2001), that states the following: let P^j $1 \leq j \leq n$, n d-dimensional independent and identically distributed random vectors and $\|\cdot\|_p$ the Minkowski norm with exponent p. If the P^j are distributed in $[0,1]^d$ then there exists a constant C_p independent of the distribution of the P^j such that

Then, there is the surprising fact that on average the

$$C_p \leq \lim_{d \rightarrow \infty} E \left(\frac{\max_j \|P^{(j)}\|_p - \min_j \|P^{(j)}\|_p}{d^{1/p}} \right) \leq (n-1) \cdot C_p.$$

$$\max_j \|P^{(j)}\|_p - \min_j \|P^{(j)}\|_p$$

contrast grows

as $d^{1/p-1/2}$. As a result, the contrast converges to a constant when the dimension increases and when the euclidean distance is used. For the L_1 norm, it increases as \sqrt{d} , for the euclidean norm ($p=2$) it remains constant and for norms with $p \geq 3$ it tends towards zero. Thus, the conclusion is that for L_p metrics with $p \geq 3$ the NN search in a high dimensional space tends to be meaningless. In other words, distance loses its discriminative power between the notions of close and far. So, on average the ratio between the contrast and $d^{1/p-1/2}$ is bounded and these bounds depend on the value of p. Furthermore, if the number of points n is large, the upper bound may be very large too. This value is much closer though to the lower bound than to the upper bound.

C. Concentration of fractional norms

Aggarwal extended Hinneburg's result to fractional p-norms (François et al., 2007; Aggarwal et al., 2001). The theorem states that if P^j $1 \leq j \leq n$ are n d-dimensional independent random vectors distributed over $[0,1]^d$ then there exists a constant C independent of p and d such that

$$C \sqrt{\frac{1}{2p+1}} \leq \lim_{d \rightarrow \infty} E \left(\frac{\max_j \|P^{(j)}\|_p - \min_j \|P^{(j)}\|_p}{\min_j \|P^{(j)}\|_p} \right) \cdot \sqrt{d} \leq (n-1) \cdot C \cdot \sqrt{\frac{1}{2p+1}}.$$

Aggarwal notes that the constant $\sqrt{1/(2p+1)}$ may play a valuable role in affecting the relative contrast and confirmed it experimentally with synthetic data sets. It was also concluded that on average fractional norms provide better contrast than Minkowski norms in terms of relative distance. Finally, Skala (2005) showed that the ratio

$$\rho_p(d) = \frac{E(\|X\|_p)^2}{2 \text{Var}(\|X\|_p)},$$

increases linearly with the dimension d . Here X is a random vector whose components are independent and identically distributed.

IV. EXPERIMENTS

Figures 3 and 4 demonstrate how the TPIE KDB-tree (Arge et al., 2002) behaves when the data set size is 20,000 and 1,000,000 points and we perform range queries that contain the number of points shown (of course with the relevant side length in each dimension).

As low as in 8 dimensions TPIE KDB-tree must visit all the created nodes in order to find the desired number of points. This result demonstrates the appearance of the dimensionality curse phenomenon, since a plain sequential scan is more efficient than using the KDB-tree. When the dataset is 1,000,000 points this phenomenon occurs when the dimensionality is 16.

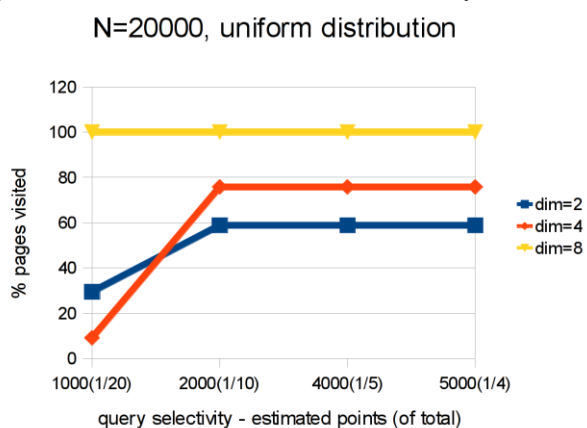


Figure 3. Percentage of visited pages for varying query selectivity and dimensionality (N=20000)

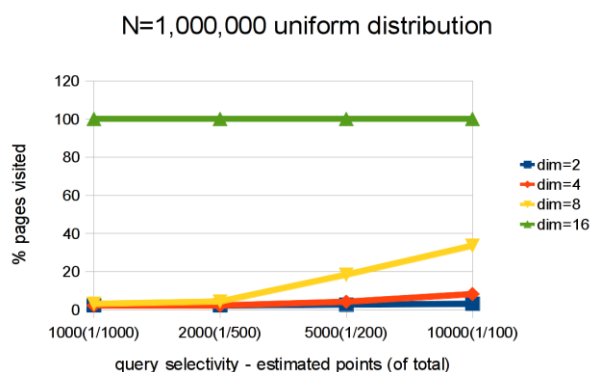


Figure 4. Percentage of visited pages for varying query selectivity and dimensionality (N=1,000,000)

V. CONCLUSIONS

In this paper we reviewed in depth the current findings on the study of high dimensional spaces. We gave many different explanations of the notion of the dimensionality curse. Finally, we demonstrated how the KDB-tree behaves in low to medium dimensions and how the dimensionality curse appears even in low dimensions and small database sizes.

REFERENCES

- [Charu C. Aggarwal: Re-designing Distance Functions and Distance-Based Applications for High Dimensional Data. SIGMOD Record 30\(1\): 13-18 \(2001a\)](#)
- [Charu C. Aggarwal, Alexander Hinneburg, Daniel A. Keim: On the Surprising Behavior of Distance Metrics in High Dimensional Spaces. ICDT 2001: 420-434](#)
- [Lars Arge, Octavian Procopiuc, Jeffrey Scott Vitter: Implementing I/O-efficient Data Structures Using TPIE. ESA 2002:88-100](#)
- [Jon Louis Bentley: Multidimensional Binary Search Trees Used for Associative Searching. Commun. ACM \(CACM\) 18\(9\):509-517 \(1975\)](#)
- [Kevin S. Beyer, Jonathan Goldstein, Raghu Ramakrishnan, Uri Shaft: When Is "Nearest Neighbor" Meaningful? ICDT 1999: 217-235](#)
- [Pierre Demartines, "Analyse de Donnees par Reseaux de Neurones Auto-Organises," PhD dissertation, Institut Nat'l Polytechnique de Grenoble, Grenoble, France, 1994 \(in French\)](#)
- [Damien Francois, Vincent Wertz, Michel Verleysen: The Concentration of Fractional Distances. IEEE Trans. Knowl. Data Eng. 19\(7\): 873-886 \(2007\)](#)
- [Michel Ledoux: The Concentration of Measure Phenomenon. American Mathematical Society 2001](#)
- [John Robinson: The K-D-B-tree: a search structure for large multidimensional dynamic indexes. Sigmod 1981](#)
- [M. Skala, "Measuring the Difficulty of Distance-Based Indexing," Proc. 12th Int'l Conf. String Processing and Information Retrieval \(SPIRE '05\), M.P. Consens and G. Navarro, eds., pp. 103-114, Nov.2005](#)
- [Roger Weber, Hans-Jörg Schek, Stephen Blott: A Quantitative Analysis and Performance Study for Similarity-Search Methods in High-Dimensional Spaces. VLDB 1998: 194-205](#)