

# Lecture Notes in Computer Science

2341

Edited by G. Goos, J. Hartmanis, and J. van Leeuwen

**Springer**

*Berlin*

*Heidelberg*

*New York*

*Barcelona*

*Hong Kong*

*London*

*Milan*

*Paris*

*Tokyo*

Cui Yu

# High-Dimensional Indexing

Transformational Approaches  
to High-Dimensional Range and Similarity Searches



Springer

## Series Editors

Gerhard Goos, Karlsruhe University, Germany  
Juris Hartmanis, Cornell University, NY, USA  
Jan van Leeuwen, Utrecht University, The Netherlands

## Author

Cui Yu  
Monmouth University, Department of Computer Science  
West Long Branch, NJ 07764, USA

National University of Singapore, Department of Computer Science  
Kent Ridge, Singapore 117543, Singapore

E-mail: [cyu@monmouth.edu](mailto:cyu@monmouth.edu)

## Cataloging-in-Publication Data applied for

Bibliographic information published by Die Deutsche Bibliothek

Die Deutsche Bibliothek lists this publication in the Deutsche Nationalbibliografie;  
detailed bibliographic data is available in the Internet at [<http://dnb.ddb.de>](http://dnb.ddb.de).

CR Subject Classification (1998): H.3.1, H.2.8, H.3, H.2, E.2, E.1, H.4, H.5.1

ISSN 0302-9743

ISBN 3-540-44199-9 Springer-Verlag Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer-Verlag. Violations are liable for prosecution under the German Copyright Law.

Springer-Verlag Berlin Heidelberg New York  
a member of BertelsmannSpringer Science+Business Media GmbH

<http://www.springer.de>

© Springer-Verlag Berlin Heidelberg 2002  
Printed in Germany

Typesetting: Camera-ready by author, data conversion by Boller Mediendesign  
Printed on acid-free paper      SPIN: 10869935      06/3142      5 4 3 2 1 0

# Preface

Many new applications, such as multimedia databases, employ the so-called feature transformation which transforms important features or properties of data objects into high-dimensional points. Searching for ‘similar’ objects based on these features is thus a search of points in this feature space. Another high-dimensional database example is stock price information systems, where time series data are stored and searched as high-dimensional data points. To support efficient query processing and knowledge discovery in these high-dimensional databases, high-dimensional indexes are required to prune the search space and efficient similarity join strategies employing these indexes have to be designed.

High-dimensional indexing has been considered an important means to facilitate fast query processing in data mining, fast retrieval, and similarity search in image and scientific databases. Existing multi-dimensional indexes such as R-trees are not scalable in terms of the number of dimensions. It has been observed that the performances of R-tree-based index structures deteriorate rapidly when the dimensionality of data is high [11, 12]. This is due to rapid growth in overlap in the directory with respect to growing dimensionality of data, requiring a large number of subtrees to be searched for each query. The problem is further exacerbated by the fact that a small high-dimensional query covering a very small fraction of the data space has actually a very large query width along each dimension. Larger query widths imply that more subtrees need to be searched. In this monograph, we study the problem of high-dimensional indexing to support range, similarity, and  $K$ -nearest neighbor (KNN) queries, and similarity joins.

To efficiently support window/range queries, we propose a simple and yet efficient transformation-based method called the  $iMinMax(\theta)$ . The method maps points in high-dimensional spaces to single dimensional values determined by their maximum or minimum values among all dimensions. With such representations, we are able to index high-dimensional data points using a conventional  $B^+$ -tree. By varying the tuning ‘knob’,  $\theta$ , we can obtain a different family of  $iMinMax$  structures that are optimized for different distributions of data sets. Hence, the method is tunable to yield best performance based on data distributions. For a  $d$ -dimensional space, a window query needs to be transformed into  $d$  subqueries. However, some of these subqueries can

be pruned away without evaluation, further enhancing the efficiency of the scheme. Extensive experiments were conducted, and experimental comparison with other existing methods such as the VA-file and Pyramid-tree provides an insight on the efficiency of the proposed method.

To efficiently support similarity or  $K$ -nearest neighbor (KNN) queries, we propose a specialized metric-based index called iDistance, and an extension of the iMinMax( $\theta$ ). In the iDistance, a metric-based index, the high-dimensional space is split into partitions, and each partition is associated with an ‘anchor’ point (called a reference point) whereby other points in the same partitions can be made reference to. With such a representation, the transformed points can then be indexed using a  $B^+$ -tree, and KNN search in the high-dimensional space is performed as a sequence of increasingly larger range queries on the single dimensional space. Such an approach supports efficient filtering of data points that are obviously not in the answer set without incurring expensive distance computation. Furthermore, it facilitates fast initial response time by providing users with approximate answers *online* that are progressively refined till all correct answers are obtained (unless the users terminate prematurely). Unlike KNN search, similarity range search on iDistance is straightforward and is performed as a spherical range query with fixed search radius. Extensive experiments were conducted, and experimental results show that the iDistance is an efficient index structure for nearest neighbor search.

The iMinMax( $\theta$ ) is designed as a generic structure for high-dimensional indexing. To extend the iMinMax( $\theta$ ) for KNN search, we design KNN processing strategies based on range search to retrieve approximate nearest neighbor data points with respect to a given query point. With proper data sampling, accuracy up to 90% can be supported very efficiently. For a more accurate retrieval, bigger search ranges must be used, which is less efficient.

In conclusion, both iMinMax( $\theta$ ) and iDistance methods are flexible, efficient, and easy to implement. Both methods can be crafted into existing DBMSs easily. This monograph shows that efficient indexes need not necessarily be complex, and the  $B^+$ -tree, which was designed for traditional single dimensional data, could be just as efficient for high-dimensional indexing. The advantage of using the  $B^+$ -tree is obvious. The  $B^+$ -tree is well tested and optimized, and so are its other related components such as concurrency control, space allocation strategies for index and leaf nodes, etc. Most importantly, it is supported by most commercial DBMSs. A note of caution is that, while it may appear to be straightforward to apply transformation on any data set to reuse  $B^+$ -trees, guaranteeing good performance is a non-trivial task. In other words, a careless choice of transformation scheme can lead to very poor performance. I hope this monograph will provide a reference for and benefit those who intend to work on high-dimensional indexing.

I am indebted to a number of people who have assisted me in one way or another in materializing this monograph. First of all, I wish to express my

appreciation to Beng Chin Ooi, for his insight, encouragement, and patience. He has taught me a great deal, instilled courage and confidence in me, and shaped my research capability. Without him, this monograph, which is an extended version of my PhD thesis [104], would not have materialized.

I would like to thank Kian-Lee Tan and Stéphane Bressan for their advice and suggestions. Kian-Lee has also proof-read this monograph and provided detailed comments that greatly improved the literary style of this monograph. I would like to thank H.V. Jagadish, for his insight, comments, and suggestions regarding iDistance; Rudolf Bayer and Mario Nascimento, for their comments and suggestions concerning the thesis; and many kind colleagues, for making their source codes available. I would like to thank Shuguang Wang, Anirban Mondal, Hengtao Shen, and Bin Cui, and the editorial staff of Springer-Verlag for their assistance in preparing this monograph. I would like to thank the School of Computing, National University of Singapore, for providing me with a graduate scholarship and facility for completing this monograph.

Last but not least, I would like to thank my family for their support, and I would like to dedicate this monograph to my parents for their love.

May 2002

*Cui Yu*

# Contents

<b>1. Introduction</b>	1
1.1 High-Dimensional Applications	1
1.2 Motivations	4
1.3 The Objectives and Contributions	7
1.4 Organization of the Monograph	8
<b>2. High-Dimensional Indexing</b>	9
2.1 Introduction	9
2.2 Hierarchical Multi-dimensional Indexes	11
2.2.1 The R-tree	11
2.2.2 Use of Larger Fanouts	14
2.2.3 Use of Bounding Spheres	15
2.2.4 The <i>kd</i> -tree	16
2.3 Dimensionality Reduction	17
2.3.1 Indexing Based on Important Attributes	18
2.3.2 Dimensionality Reduction Based on Clustering	18
2.3.3 Mapping from Higher to Lower Dimension	20
2.3.4 Indexing Based on Single Attribute Values	22
2.4 Filtering and Refining	26
2.4.1 Multi-step Processing	26
2.4.2 Quantization	27
2.5 Indexing Based on Metric Distance	29
2.6 Approximate Nearest Neighbor Search	32
2.7 Summary	33
<b>3. Indexing the Edges – A Simple and Yet Efficient Approach to High-Dimensional Range Search</b>	37
3.1 Introduction	37
3.2 Basic Concept of iMinMax	38
3.2.1 Sequential Scan	41
3.2.2 Indexing Based on Max/Min	41
3.2.3 Indexing Based on iMax	42
3.2.4 Preliminary Empirical Study	45
3.3 The iMinMax Method	46



3.4	Indexing Based on iMinMax	47
3.5	The iMinMax( $\theta$ )	49
3.6	Processing of Range Queries	52
3.7	iMinMax( $\theta$ ) Search Algorithms	57
3.7.1	Point Search Algorithm	57
3.7.2	Range Search Algorithm	57
3.7.3	Discussion on Update Algorithms	58
3.8	The iMinMax( $\theta_i$ )	58
3.8.1	Determining $\theta_i$	59
3.8.2	Refining $\theta_i$	62
3.8.3	Generating the Index Key	63
3.9	Summary	64
<b>4.</b>	<b>Performance Study of Window Queries</b>	<b>65</b>
4.1	Introduction	65
4.2	Implementation	65
4.3	Generation of Data Sets and Window Queries	66
4.4	Experiment Setup	66
4.5	Effect of the Number of Dimensions	67
4.6	Effect of Data Size	69
4.7	Effect of Skewed Data Distributions	70
4.8	Effect of Buffer Space	76
4.9	CPU Cost	77
4.10	Effect of $\theta_i$	78
4.11	Effect of Quantization on Feature Vectors	80
4.12	Summary	83
<b>5.</b>	<b>Indexing the Relative Distance – An Efficient Approach to KNN Search</b>	<b>85</b>
5.1	Introduction	85
5.2	Background and Notations	86
5.3	The iDistance	87
5.3.1	The Big Picture	88
5.3.2	The Data Structure	90
5.3.3	KNN Search in iDistance	91
5.4	Selection of Reference Points and Data Space Partitioning	95
5.4.1	Space-Based Partitioning	96
5.4.2	Data-Based Partitioning	99
5.5	Exploiting iDistance in Similarity Joins	102
5.5.1	Join Strategies	102
5.5.2	Similarity Join Strategies Based on iDistance	103
5.6	Summary	107

<b>6. Similarity Range and Approximate KNN Searches with iMinMax</b>	109
6.1 Introduction	109
6.2 A Quick Review of iMinMax( $\theta$ )	109
6.3 Approximate KNN Processing with iMinMax	110
6.4 Quality of KNN Answers Using iMinMax	115
6.4.1 Accuracy of KNN Search	118
6.4.2 Bounding Box Vs. Bounding Sphere	118
6.4.3 Effect of Search Radius	118
6.5 Summary	120
<b>7. Performance Study of Similarity Queries</b>	123
7.1 Introduction	123
7.2 Experiment Setup	123
7.3 Effect of Search Radius on Query Accuracy	123
7.4 Effect of Reference Points on Space-Based Partitioning Schemes	126
7.5 Effect of Reference Points on Cluster-Based Partitioning Schemes	127
7.6 CPU Cost	131
7.7 Comparative Study of iDistance and iMinMax	133
7.8 Comparative Study of iDistance and A-tree	134
7.9 Comparative Study of the iDistance and M-tree	136
7.10 iDistance – A Good Candidate for Main Memory Indexing?	137
7.11 Summary	139
<b>8. Conclusions</b>	141
8.1 Contributions	141
8.2 Single-Dimensional Attribute Value Based Indexing	141
8.3 Metric-Based Indexing	142
8.4 Discussion on Future Work	143
<b>References</b>	145