Claudia Bauzer Medeiros
Max Egenhofer
Elisa Bertino (Eds.)

# Advances in Spatial and Temporal Databases

**9th International Symposium, SSTD 2005**
**Angra dos Reis, Brazil, August 2005**
**Proceedings**

Springer

Claudia Bauzer Medeiros    Max Egenhofer
Elisa Bertino (Eds.)

# Advances in Spatial and Temporal Databases

9th International Symposium, SSTD 2005
Angra dos Reis, Brazil, August 22-24, 2005
Proceedings

Springer

Volume Editors

Claudia Bauzer Medeiros
University of Campinas, Institute of Computing
CP 6176, 13084-971 Campinas, Brazil
E-mail: cmbm@ic.unicamp.br

Max Egenhofer
University of Maine
National Center for Geographic Information and Analysis
348 Boardman Hall, Orono, ME 04469-5711, USA
E-mail: max@spatial.maine.edu

Elisa Bertino
Purdue University, Department of Computer Science
West Lafayette, IN, USA
E-mail: bertino@cerias.purdue.edu

# Lecture Notes in Computer Science 3633

*Commenced Publication in 1973*
Founding and Former Series Editors:
Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen

## Editorial Board

# Preface

It is our great pleasure to introduce the papers of the proceedings of the 9th International Symposium on Spatial and Temporal Databases – SSTD 2005. This year's symposium continues the tradition of being the premier forum for the presentation of research results and experience reports on leading edge issues of spatial and temporal database systems, including data models, systems, applications and theory. The mission of the symposium is to share innovative solutions that fulfill the needs of novel applications and heterogeneous environments and identify new directions for future research and development. SSTD 2005 gives researchers and practitioners a unique opportunity to share their perspectives with others interested in the various aspects of database systems for managing spatial and temporal data and for supporting their applications.

A total of 77 papers were submitted this year from several countries. After a thorough review process, the program committee accepted 24 papers covering a variety of topics, including indexing techniques and query processing, mobile environments and moving objects, and spatial and temporal data streams. We are very pleased with the variety of the symposium's topics, and we are proud of the resulting strong program.

Many people contributed to the success of the SSTD 2005 program. First of all, we would like to thank the authors for providing the content of the program, and all the members of the program committee and the additional reviewers, for their detailed comments. Philippe Rigaux was of help in adding functions to his program MyReview, which was used in the reviewing process. We would also like to express our gratitude to Gilberto Câmara, the general chair of SSTD 2005, for his constant guidance and advice on many organizational aspects of the symposium and for his work on the local arrangements. Finally, we would like to thank our sponsors (notably INPE – the Brazilian National Institute for Space Research) who have enabled us to hold a successful meeting. We are also grateful for the support of the Brazilian Computer Society (SBC).

We hope that you find this program to be both beneficial and enjoyable and that the symposium provides you with the opportunity to meet other researchers and practitioners from institutions around the world. Enjoy!!

August 2005
Claudia Bauzer Medeiros,
Max Egenhofer,
Elisa Bertino

# Organization

SSTD 2005 was organized by the Department of Image Processing of INPE – the National Institute for Space Research (Instituto Nacional de Pesquisas Espaciais), Brazil.

## Executive Committee

General Conference Chair:    Gilberto Câmara (Department of Image Processing, INPE, Brazil)

Program Chairs:    Elisa Bertino (Department of Computer Sciences, Purdue University, USA)

Max Egenhofer (NCGIA, University of Maine, USA)

Claudia Bauzer Medeiros (Institute of Computing, University of Campinas, Brazil)

## Program Committee

| | |
|---|---|
| Amr El Abbadi | UC Santa Barbara, USA |
| Walid G. Aref | Purdue University, USA |
| Alessandro Artale | University of Bolzano, Italy |
| Alberto Belussi | University of Verona, Italy |
| Michela Bertolotto | University College Dublin, Ireland |
| Gilberto Camara | National Institute for Space Research, Brazil |
| Marco Casanova | Dept Informatics, PUC-Rio, Brazil |
| Barbara Catania | University of Genoa, Italy |
| Christophe Claramunt | Naval Academy Research Institute, France |
| Matt Duckham | University of Melbourne, Australia |
| Fred Fonseca | Penn State University, USA |
| Fosca Giannotti | CNR, Italy |
| Ralf Hartmut Güting | University of Hagen, Germany |
| Kathleen Hornsby | University of Maine, USA |
| Christian S. Jensen | Aalborg University, Denmark |
| Christopher Jones | Cardiff University, UK |
| Daniel Keim | University of Constance, Germany |
| Eamonn Keogh | University of California, Riverside, USA |
| George Kollios | Boston University, USA |
| Bart Kuijpers | University of Limburg, Belgium |
| Mario Nascimento | University of Alberta, Canada |
| Raymond Ng | University of British Columbia, Canada |

| | |
|---|---|
| Silvia Nittel | NCGIA, University of Maine, USA |
| Beng Chin Ooi | National University of Singapore, Singapore |
| Peter van Oosterom | Delft Univ. of Technology, The Netherlands |
| Dimitris Papadias | Hong Kong Univ. of Science and Technology, Hong Kong |
| Jignesh Patel | University of Michigan, USA |
| Sunil Prabhakar | Purdue University, USA |
| Philippe Rigaux | University Paris IX, France |
| Andrea Rodríguez-Tastets | University of Concepción, Chile |
| Ana Carolina Salgado | University of Pernambuco, Brazil |
| George Samaras | University of Cyprus, Cyprus |
| Peter Scheuermann | Northwestern University, USA |
| Markus Schneider | University of Florida, USA |
| Bernhard Seeger | University of Marburg, Germany |
| Cyrus Shahabi | University of Southern California, USA |
| Shashi Shekhar | University of Minnesota, Minneapolis, USA |
| Rick Snodgrass | University of Arizona, USA |
| Stefano Spaccapietra | EPFL Lausanne, Switzerland |
| Paolo Terenziani | Università del Piemonte Orientale, Italy |
| Yannis Theodoridis | University of Piraeus, Greece |
| Nectaria Tryfona | Computer Technology Institute, Greece |
| Michalis Vazirgiannis | Athens Univ. of Economics and Business, Greece |
| Agnes Voisard | Fraunhofer ISST and FU Berlin, Germany |
| Ouri Wolfson | University of Illinois at Chicago, USA |

## Additional Referees

| | | |
|---|---|---|
| Nagender Bandi | Feifei Li | Joern Schneidewind |
| Benjamin Bustos | Xiang Lian | Chengyu Sun |
| Hu Cao | Dan Lin | Valeria C. Times |
| Reynold Cheng | Juhong Liu | Goce Trajcevski |
| Alminas Civilis | Andrei Lopatenko | Yicheng Tu |
| Stephen Cole | Anna Maddalena | Tao Wan |
| Carlo Combi | Florian Mansmann | David Yang |
| Stephane Coulondre | Ahmed Metwally | Kiyoung Yang |
| Pier Luigi Dragotti | Mohamed Mokbel | Huabei Yin |
| Ying Feng | Kyriakos Mouratidis | Man Lung Yiu |
| Robson N. Fidalgo | Guillaume Noel | Yuni Xia |
| Elias Frentzos | Andrea Nucita | Xiaopeng Xiong |
| Marios Hadjieleftheriou | Nikos Pelekis | Bo Xu |
| Christoph Heinz | Paola Podesta | Mingwu Zhang |
| Xuegang Huang | Gianna Reggio | Hartmut Ziegler |
| Mohammad R. Kolahdouzan | Mehdi Sharifzadeh | |
| | Sarvjeet Singh | |

# Sponsoring Institutions

Sponsor – National Institute of Space Research – INPE – Brazil
    Support – Brazilian Computer Society (SBC)

# Lecture Notes in Computer Science

For information about Vols. 1–3511

please contact your bookseller or Springer

# Table of Contents

## Query Optimization and Simulation

## Advanced Query Processing I

## Spatial/Temporal Data Streams

## Advanced Query Processing II

## Indexing Schemes and Structures

## Novel Applications and Real Systems

## Moving Objects and Mobile Environments

## Advanced Query Processing III

# Selectivity Estimation of High Dimensional Window Queries via Clustering

Christian Böhm, Hans-Peter Kriegel, Peer Kröger, and Petra Linhart

Institute for Computer Science, University of Munich, Germany
{boehm, kriegel, kroegerp, linhart}@dbs.ifi.lmu.de

**Abstract.** Query optimization is an important functionality of modern database systems and often based on estimating the selectivity of queries before actually executing them. Well-known techniques for estimating the result set size of a query are sampling and histogram-based solutions. Sampling-based approaches heavily depend on the size of the drawn sample which causes a trade-off between the quality of the estimation and the time in which the estimation can be executed for large data sets. Histogram-based techniques eliminate this problem but are limited to low-dimensional data sets. They either assume that all attributes are independent which is rarely true for real-world data or else get very inefficient for high-dimensional data. In this paper we present the first multivariate parametric method for estimating the selectivity of window queries for large and high-dimensional data sets. We use clustering to compress the data by generating a precise model of the data using multivariate Gaussian distributions. Additionally, we show efficient techniques to evaluate a window query against the Gaussian distributions we generated. Our experimental evaluation shows that this approach is significantly more efficient for multidimensional data than all previous approaches.

## 1 Introduction

The storage and management of vectors of a multidimensional feature space has become an important basic functionality of a database system. Advanced applications such as multimedia [1], CAD [2], molecular biology [3], etc. require efficient and effective methods for content based similarity search and data mining. Such methods are typically based on feature vectors of moderate or high dimensionality. Although a vast number of index structures [4,5] and access methods [6] for vector data has been proposed, database management systems do not yet support the storage and retrieval of vector data in the same way as relational data from applications such as accounting and billing. In order to give full support to advanced applications the database system needs efficient and effective techniques for query optimization. One of the most important challenges in query optimization is the estimation of the selectivity of a query predicate. While a number of techniques to model the data distribution and thus to estimate the selectivity are known for one- and low-dimensional data spaces, this is still an unsolved problem for data spaces of medium to high dimensionality.

Three different paradigms of data modelling for selectivity estimation in general can be distinguished: Histograms, sampling, and parametric techniques. Of those three, only sampling can be directly applied without modification in higher dimensional data spaces. Many different sampling methods have been proposed. They share the common idea to evaluate the predicate on top of a small subset of the actual database objects and to extrapolate the observed selectivity. The well-known techniques differ in the way how the sample is drawn as well as in the determination of the suitable size of the sample. The general drawback of sampling techniques is that the accuracy of the result is strictly limited by the sample rate. To get an accurate estimation of the selectivity, a large sample (>10%) of the database is required. To evaluate the query on top of the large sample is not much cheaper than to evaluate it on the original data set which limits its usefulness for query optimization.

Histogram techniques, the most prevalent paradigm to model the data distribution in the one-dimensional case, have a different problem. This concept is very difficult to be carried over to the multidimensional case, even for low or moderate dimensional data. One way to adapt one-dimensional histograms to multidimensional data is to describe the distribution of the individual attributes of the vectors independently by usual histograms. These histograms are sometimes called marginal distributions. In this case, the selectivity of multidimensional queries can be determined easily provided that the attributes are statistically independent, i.e. neither correlated nor clustered. Real-world data sets, however, rarely fulfill this condition. Another approach is to partition the data space by a multidimensional grid and to assign a histogram bin to each grid cell. This approach may be possible for two- and three-dimensional spaces. However, for higher dimensional data this method becomes inefficient and ineffective since the number of grid cells is exponential in the dimensionality. Techniques of dimensionality reduction such as Fourier transformation, wavelets, principal component analysis or space-filling curves (Z-ordering, Hilbert) may reduce this problem to some extent. The possible problem reduction, however, is limited by the intrinsic dimensionality of the data set.

The idea of parametric techniques is to describe the data distribution by curves (functions) which have been fitted into the data set. In most cases Gaussian functions (normal distributions) are used. Instead of using one single Gaussian, a set of multivariate Gaussians can be fitted into the data set which makes the technique more accurate. Each Gaussian is then described by three parameters (mean, variance and the relative weight of the Gaussian in the ensemble). This approach can be transferred into the multidimensional case by two techniques. Like described above for histograms, the marginal distribution of each attribute can be modelled independently by a set of Gaussians. The multidimensional query selectivity can be estimated by combining the marginal distributions. This approach leads to similar problems like marginal histograms.

Therefore, our solution is different. Our technique directly models the multidimensional data distribution by a set of multivariate Gaussian functions. There are two options to use the Gaussian primitives: The Gaussians can either be used

with a matrix containing both variances and covariances or with a vector of the multivariate variances only. As we will discuss later, both approaches have their advantages and disadvantages. When using Gaussians with covariance matrix, the data distribution can be described more accurately by a single primitive. On the other side, more storage is needed for the covariance matrices ($O(d^2)$ for each Gaussian) compared to the variance vector approach ($O(d)$ for each Gaussian). Moreover, the processing cost for reading the parameters and for the determination of the estimated selectivity is much higher when covariance matrices are used. Let us note that, unlike the approaches using marginal distributions, our Gaussian technique with no covariance matrix does not rely on the attribute independence assumption. This technique assumes attribute independence for each individual Gaussian primitive only, but places no constraints to the overall data distribution. We will discuss this issue in more detail in Section 4, an experimental validation is given in Section 5.
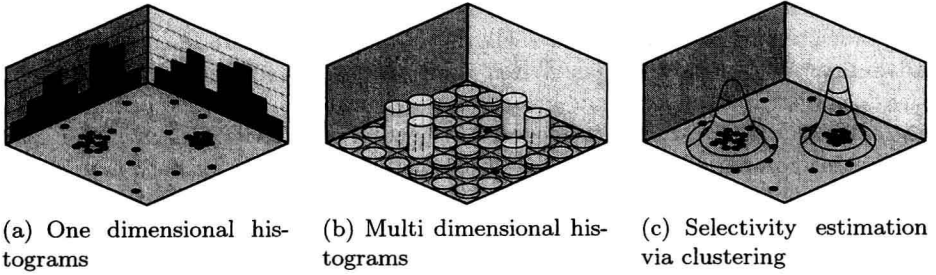
To obtain a collection of Gaussians distributions we apply a clustering algorithm. Clustering is the task of grouping vectors into different subsets (the clusters) such that the intra-cluster similarity is maximized and the inter-cluster similarity is minimized. That means points belonging to the same cluster are close together whereas points of different clusters are far away from each other. Many different algorithms have been proposed such as k-means [7], single-link [8], density-based clustering [9,10] and many others. Most of these algorithms use a point as a representative of each cluster. In contrast, the EM clustering algorithm (expectation maximization) [11] uses a multivariate Gaussian function as a cluster representative. We will discuss the suitability of different variants of the EM algorithm for our problem of getting a good approximation of the actual data distribution.

To summarize our contribution, we propose in this paper a new cost model for estimating the selectivity of multidimensional queries on top of vector data of medium to high dimensionality. The data distribution is represented by a set of multivariate Gaussian functions that have been determined using the EM clustering algorithm. We develop two methods for estimating the selectivity of window queries and range queries using the multivariate Gaussians. We demonstrate experimentally the superiority of our approach over competitive cost models based on histograms and sampling. The remainder of our paper is organized as follows: In Section 2 we discuss related work on selectivity estimation and point out our contribution. Section 3 and 4 describes in detail our proposed methods to find a representation of the data distribution by an ensemble of multivariate Gaussian functions using EM clustering and to estimate the selectivity on top of this model. Section 5 contains the experimental evaluation, and section 6 concludes our paper.

## 2   Related Work

In this chapter, we review current approaches for selectivity estimation and discuss their potentials.

(a) One dimensional his-    (b) Multi dimensional his-    (c) Selectivity  estimation
tograms                      tograms                        via clustering

**Fig. 1.** Visualization of different concepts for selectivity estimation

## 2.1   Review

Recent work on selectivity estimation can be categorized into three classes,
namely histogram-based methods, sampling-based methods, and parametric
methods. In the following, we review and discuss the most important repre-
sentatives of each class briefly.

**Histogram-based Methods.** The most widespread approach for selectivity
estimation in practice is the use of histograms. In general, the data space is par-
titioned into buckets, and the frequency of points inside each bucket is computed.
We can distinguish between one-dimensional and multi-dimensional histograms.

Selectivity estimation using one-dimensional histograms is based on the as-
sumption that the attributes of the data set are independent, i.e. there is no
correlation between different dimensions of the feature space. For each dimen-
sion, a histogram is built and the selectivity of a window query $q$ is estimated
in each dimension separately. The selectivity of $q$ in the full-dimensional space
is evaluated by multiplying the selectivity estimations for each attribute. Equi-
width histograms [12] compute buckets of fixed size and variable point frequency,
whereas equi-depth histograms [13] compute buckets of variable size and fixed
point frequency.

With growing dimensionality of the feature space, the recombination of one-
dimensional buckets becomes costly. Thus, in recent years, multi-dimensional
histograms have been investigated. Multi-dimensional equi-depth histograms [14]
partition the feature space into multi-dimensional buckets with variable size and
fixed point frequency. In [14] an algorithm to construct multi-dimensional equi-
depth histograms is presented that iteratively partitions the data space along
each attribute into a fixed number of buckets, where the order of attributes
is fixed. The selectivity of a window query $q$ is estimated analogously to one-
dimensional histograms taking the buckets into account that intersect with $q$.
The algorithm MHIST [15] partitions the data space along the single attributes
in a similar way, but decides in each step which attribute is partitioned rather
than processing the attributes in a fixed order.

STHoles [16] is a recent approach that proposes hierarchically organized
multi-dimensional histograms. A histogram may contain another histogram com-