

ProUD: Probabilistic Ranking in Uncertain Databases

Matthias Renz,

Thomas Bernecker and Hans-Peter Kriegel

Ludwig-Maximilians-Universität München
Munich, Germany

www.dbs.ifi.lmu.de



- **Introduction**
- Ranking on Uncertain Data
- Probabilistic Ranking Algorithm
- Experimental Evaluation
- Summary



Introduction

- Ranking Queries
 - are very important for similarity search
 - give the most relevant answers first
 - are more flexible than distance range and k NN queries
- Applications:
 - person identification
 - similarity search in multimedia databases
 - ...
- ➔ often uncertainty in data



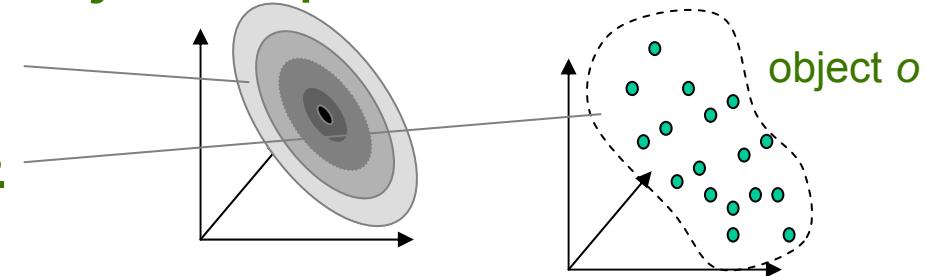
Introduction

- uncertain data caused by
 - continuously changing data
 - concurring data descriptions
 - measurement inaccuracies
 - uncertain data leads to uncertain query results
- probabilistic queries
(e.g. probabilistic ranking)
- results are associated with confidence values

- Introduction
- **Ranking on Uncertain Data**
- Probabilistic Ranking Algorithm
- Experimental Evaluation
- Summary

Ranking on Uncertain Data

- Positionally Uncertain Data
 - vector data in d -dimensional space \mathbb{R}^d
 - no unique position in \mathbb{R}^d
 - objects are represented by
 - multiple d -dimensional vectors
 - that are mutually exclusive
 - a confidence value is assigned to each vector
 - types of uncertain object representations
 - pdf (continuous)
 - discrete samples



Ranking on Uncertain Data

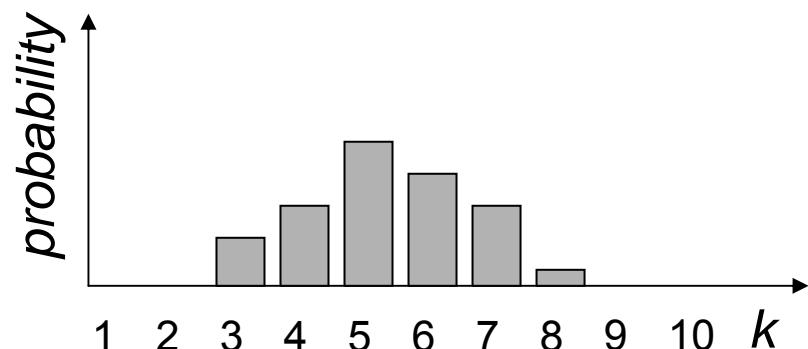
- Distance Computation for Uncertain Objects
 - two uncertain objects
$$o_i = \{o_{i,1}, \dots, o_{i,m}\} \text{ and } o_j = \{o_{j,1}, \dots, o_{j,m}\}$$
 - distance between o_i and o_j :
$$d_{uncertain}(o_i, o_j) = \{dist(o_{i,m}, o_{j,n}) \mid 1 \leq m \leq M, 1 \leq n \leq M\}$$
 - probability that the distance between o_i and o_j is less than $\varepsilon \in \mathbb{R}_0^+$:

$$P(d_{uncertain}(o_i, o_j) \leq \varepsilon) = \frac{|\{d \in d_{uncertain}(o_i, o_j) : d \leq \varepsilon\}|}{|d_{uncertain}(o_i, o_j)|}$$



Ranking on Uncertain Data

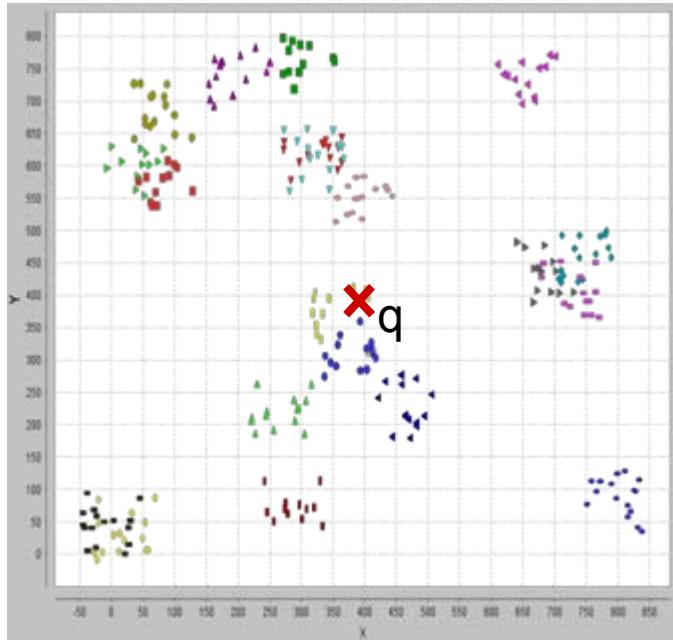
- Probabilistic Ranking
 - function $prob_ranked_q: \mathcal{D} \times \{1, \dots, N\} \rightarrow [0..1]$
 - $prob_ranked_q(o, k)$ reports the probability that object o is the k^{th} -nearest-neighbor of the query point q
- for each object:



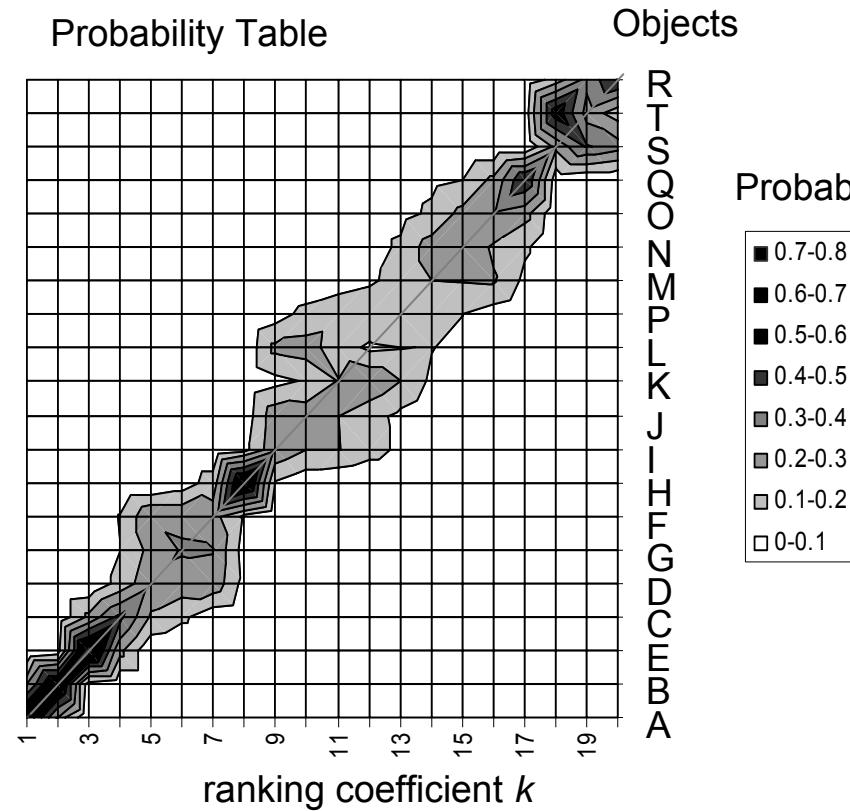


Ranking on Uncertain Data

- Probabilistic Ranking Output:
 - Example:



vector space

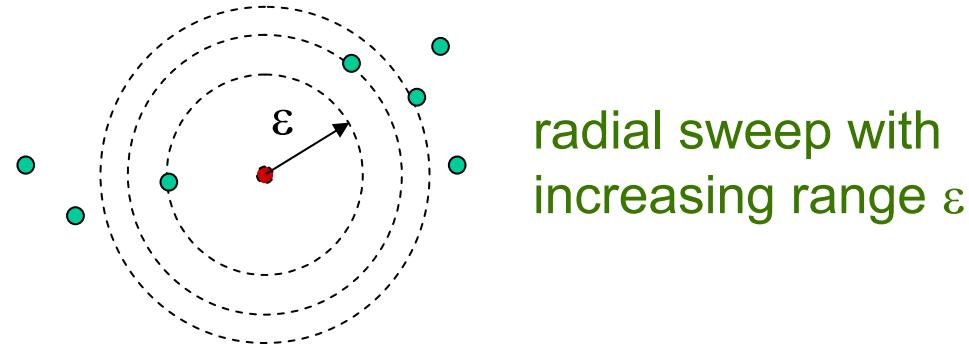


probabilistic ranking output

Outline

- Introduction
- Ranking on Uncertain Data
- **Probabilistic Ranking Algorithm**
- Experimental Evaluation
- Summary

- Iterative Probability Computation
 - ranking query on the object samples



- during the radial sweep: maintain for each object o the probability $P_o = P(d_{uncertain}(o, q) \leq \varepsilon)$
- for each accessed sample $o_{i,j}$, compute the probability $P(o_{i,j}, k)$ that exactly $(k-1)$ objects $o \neq o_i$ are within the ε -range, for $k = 1..N$.

- computation of $P(o_{i,j}, k)$:

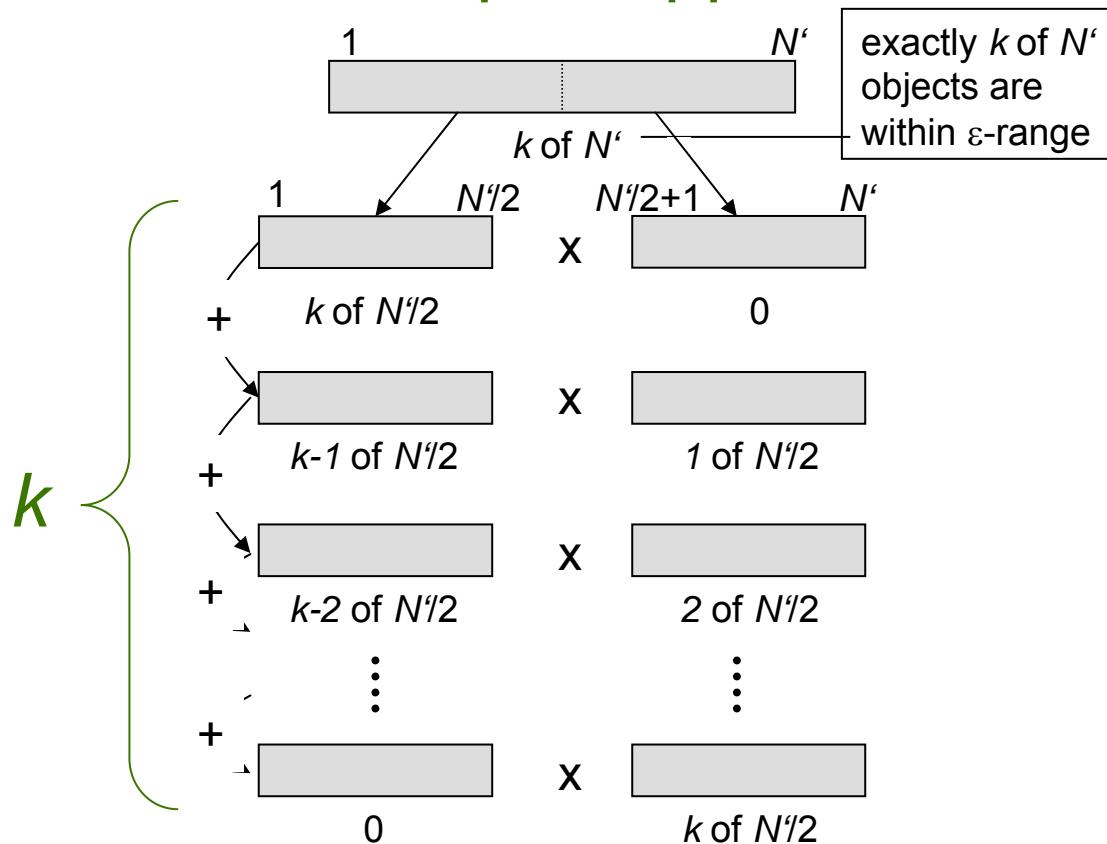
$$P(o_{i,j}, k) = \sum_{\sigma_k(i) \in S_k} \prod_{l=1..N}^{\substack{l \neq i}} \begin{cases} P_{o_l} & \text{if } o_l \in \sigma_k(i) \\ (1 - P_{o_l}) & \text{else} \end{cases}$$

- problem: a lot of possibilities to select the set $\sigma_k(i)$ of $(k-1)$ objects out of N objects
- reduce $N \rightarrow N'$ by pruning objects that are beyond ε



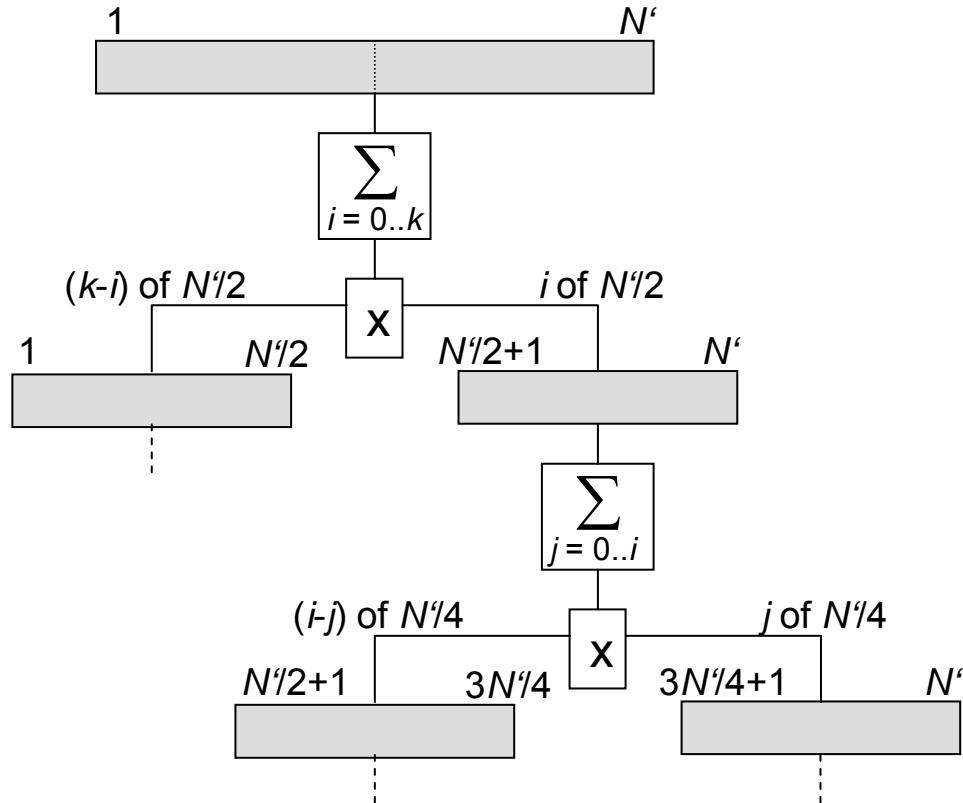
Probabilistic Ranking Algorithm

- Accelerated Probability Computation
 - divide and conquer approach





- Accelerated Probability Computation
 - divide and conquer approach

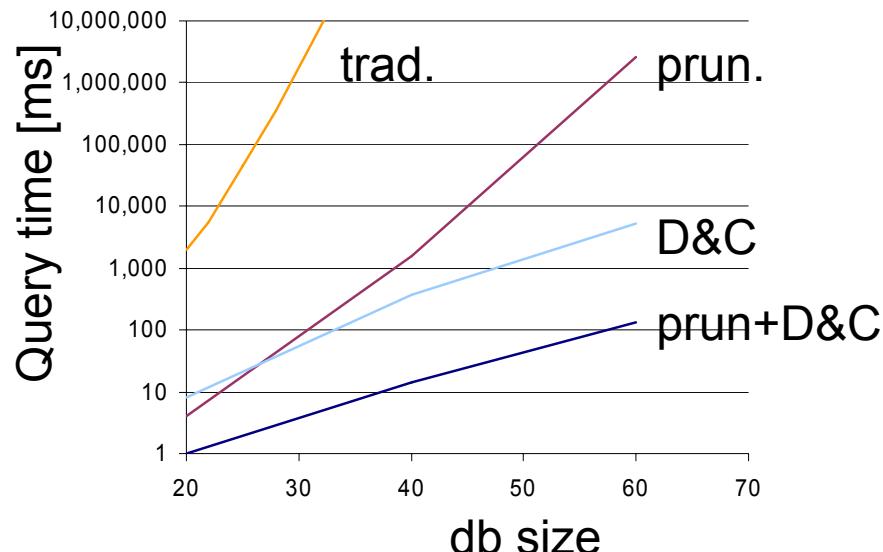


Outline

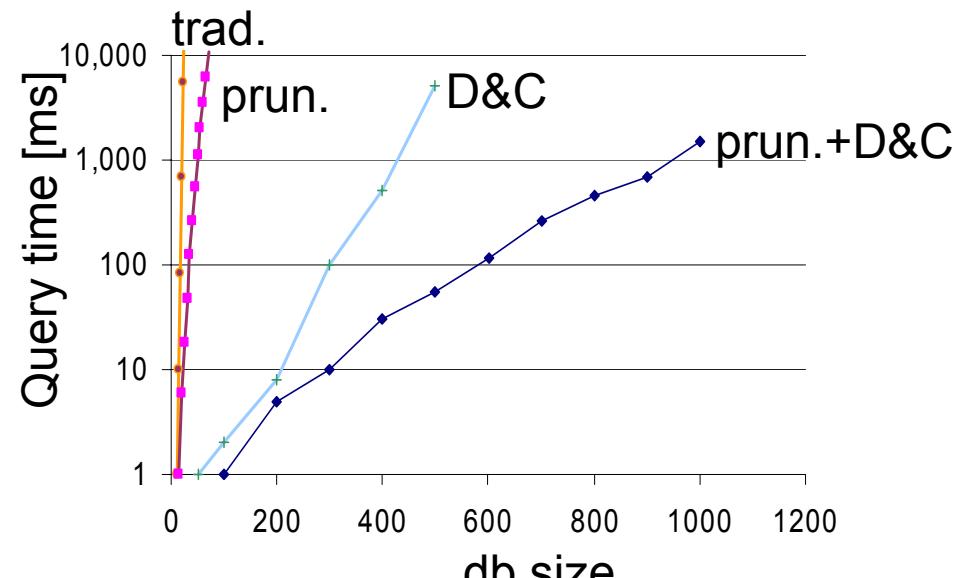
- Introduction
- Ranking on Uncertain Data
- Probabilistic Ranking Algorithm
- **Experimental Evaluation**
- Summary

Experimental Evaluation

- Performance w.r.t. database size



higher degree of uncertainty

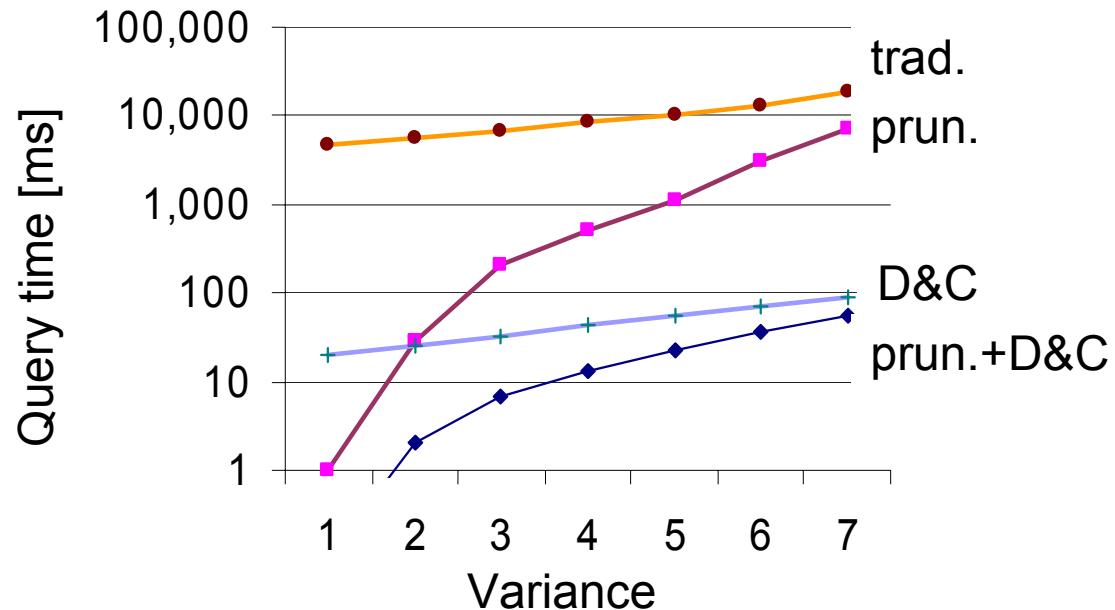


lower degree of uncertainty

- performance depends on the degree of uncertainty (object variance)

Experimental Evaluation

- Performance w.r.t. degree of uncertainty



- speed-up by pruning degenerates for higher object variances (object uncertainty)

Outline

- Introduction
- Ranking on Uncertain Data
- Probabilistic Ranking Algorithm
- Experimental Evaluation
- **Summary**



Summary

- approach to accelerate probabilistic ranking queries
- very high speed-up factor (2 orders of magnitude)
- Future Work:
 - comparison to the latest top- k query approaches



*Thank you,
for your attention,
any questions?*