# Non-Standard-Datenbanken und Data Mining

Probabilistic Spatio-Temporal Databases and Streams

Prof. Dr. Ralf Möller Universität zu Lübeck Institut für Informationssysteme



**IM FOCUS DAS LEBEN** 

# Übersicht

- Semistrukturierte Datenbanken (JSON, XML) und Volltextsuche
- Information Retrieval
- Mehrdimensionale Indexstrukturen
- Cluster-Bildung
- Einbettungstechniken
- First-n-, Top-k-, und Skyline-Anfragen
- Probabilistische Datenbanken, Anfragebeantwortung, Top-k-Anfragen und Open-World-Annahme
- Probabilistische Modellierung, Bayes-Netze, Anfragebeantwortungsalgorithmen, Lernverfahren,
- Temporale Datenbanken und das relationale Modell,
- Probabilistische Temporale Datenbanken
- SQL: neue Entwicklungen (z.B. JSON-Strukturen und Arrays), Zeitreihen (z.B. TimeScaleDB)
- Stromdatenbanken, Prinzipien der Fenster-orientierten inkrementellen Verarbeitung
- Approximationstechniken für Stromdatenverarbeitung, Stream-Mining
- Probabilistische raum-zeitliche Datenbanken und Stromdatenverarbeitungsssysteme: Anfragen und Indexstrukturen, Raum-zeitliches Data Mining
- Von NoSQL- zu NewSQL-Datenbanken, CAP-Theorem, CALM-Theorem
- Blockchain-Datenbanken
- Analyse von Graphdaten



#### Presentation slides are largely taken from Location-aware Query Processing and Optimization: A Tutorial

Mohamed F. Mokbel

Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN, USA <u>mokbel@cs.umn.edu</u>

Walid G. Aref Department of Computer Science, Purdue University, West Lafayette, Indiana, USA. <u>aref@cs.purdue.edu</u>

Some slides (indicated) were produced by George Kollios

#### Slides have been modified or extended. Faults are mine!



# Spatio-Temporal Objects

- Moving points (extent does not matter)
  - Each object is modeled as a point (e.g., moving vehicles in a GIS based transportation system)
- Moving regions (extent matters)
  - Each object is represented by an MBR, the MBR can change as the object moves (e.g., thunderstorm, noise)



#### Location-aware Query Processing and Optimization: A Tutorial, Mohamed F. Mokbel, Walid G. Aref

# Location-aware Queries

#### Continuously report the number of cars on freeway 71-75

- Type: Range query
- Time: Present
- Duration: Continuous

#### What are my nearest McDonalds for the next hour?

- Type: Nearest-neighbor query
- Time: Future
- Duration: Continuous / Snapshot

#### Send E-coupons to all cars that I am their nearest gas station

- Type: Reverse NN query
- *Time:* **Present**
- Duration: Snapshot

What was the closest distance between Taxi A & me yesterday?

- Type: Closest-point query
- *Time:* Past

NIVERSITÄT ZU LÜBECK Institut für informationssysteme

Duration: Snapshot

Query: Moving (reference rectangle)

5

IM FOCUS DAS LEBEN

- *Objects:* Stationary (McDonalds)
- Query: Stationary (gas station)
- Objects: Moving

Query: Stationary

**Objects:** Moving

- Query: Moving
- Objects: Moving

#### **Snapshot Querying the Past**

- Examples:
  - **Temporal** Dimension: What was the location of a certain object from 7:00 AM to 10:00 AM yesterday?
  - **Spatial** Dimension: *Find all objects that were in a certain area at 7:00 AM yesterday*
  - **Spatio-temporal** Dimension: *Find all objects that were close to each other from 7:00 AM to 8:00 AM yesterday*
- Features:
  - Large number of historical trajectories
  - Persistent read-only data
  - Query spatial and/or temporal dimensions





Historical trajectories are represented by their three-dimensional Minimum Bounding Rectangle (MBR)

Time 3D R-tree can be used to index MBRs Technique simple and easy to implement Does not scale well Does not provide efficient query support for snapshot queries (aka timestamp queries)



### 3D R-Tree





Objects are somewhere in the gray rectangular regions.

#### IM FOCUS DAS LEBEN 8

# **Modeling Evolution: Historical R-Trees**





## Multi-Version Index Structures (MVR-Trees)

- Maintain an R-tree for each time instance (aka historical r-tree, HR-tree)
- R-tree nodes that are not changed across consecutive time instances are linked together (remove redundancies: MVR-tree)



• A multi-version R-tree can be combined with a 3D-R-tree to support interval queries (combination is called MV3R-Tree)



Yufei Tao and Dimitris Papadias. MV3R-Tree: A Spatio-temporal Access Method for Timestamp and Interval Queries. In Proc. VLDB-01, pp. 431-440, **2001** 

IM FOCUS DAS LEBEN 10

Location-aware Query Processing and Optimization: A Tutorial, Mohamed F. Mokbel, Walid G. Aref

# Historical R-trees (HR-trees)

An R-tree is maintained for each timestamp in history.



© George Kollios

11

# **Historical R-trees**

An R-tree is maintained for each timestamp in history.

Trees at consecutive timestamps may share branches to save space.



<sup>©</sup> George Kollios

#### Building a 3D R-tree on the Leaves of the MVR-tree

- Size of the 3D R-tree is much smaller than a complete 3D R-tree as the number of leaf nodes is significantly lower than the number of actual objects.
- Long interval queries can be processed with auxiliary 3D R-trees



# Rectangles

Problem of indexing any type of moving objects can be reduced to indexing discrete rectangles





# Optimization

- If N objects move with linear functions of time:
- Minimize total volume by splitting in equidistant points
- Given K splits you can decide the best splits in O(K log N) time.



Yufei Tao and Dimitris Papadias. MV3R-Tree: A Spatio-temporal Access Method for Timestamp and Interval Queries. In Proc. VLDB-01, pp. 431-440, **2001** 

#### Querying the Present

- Time is always NOW
- Example Queries:
  - Find the number of objects in a certain area
  - What is the current location of a certain object?
- Features:
  - Continuously changing data
  - Real-time query support is required
  - Index structures should be update-tolerant
- Present data is always accessed through continuous queries





### **Updating Index Structures**

- Traditional R-tree updates are top-down
- Updates translated to delete and insert transactions
- To support frequent updates:
  - Updates can be managed
    "inline" without the need for deletion or insertions
  - Bottom-up approaches through auxiliary index structures to locate the object identifier







IM FOCUS DAS LEBEN 17

# Querying the Future

- Examples:
  - What will my nearest restaurant be after 30 minutes?
  - Does my path conflict with any other cars for the next hour?
- Features:
  - Predict the movement through a velocity vector
  - Prediction could be valid for only a limited time horizon in the future





Location prediction seems to be a simple task in some cases:



Jonas Lüthke. Location Prediction Based on Mobility Patterns in Location Histories. Master thesis, TU Hamburg-Harburg, **2013** 

https://www.ifis.uni-luebeck.de/~moeller/publist-sts-pw-and-m/source/papers/2013/luethke13.pdf



IM FOCUS DAS LEBEN 19

Master Thesis Jonas Lüthke, TUHH, 2013

### **Location Prediction - Approach**

Location prediction seems to be a simple task in some cases:



Previous observations can enable an educated guess



### **Example: Location History Data**



#### Cabspotting data set:

- GPS coordinates collected from 563 cabs in San Francisco over 30 days
- Interval between measurements
  < 60seconds</li>
- Ten taxis selected for testing (with regard to measurement density, measurement errors)

- Spatial probability distribution could be estimated from this (e.g., GMM)
- Spatiotemporal probability distribution is needed



Embed location time series in 2*m*-dimensional space using a delay *v*:

- Time series is iteratively sampled using delay time v
- Every *m* subsequent locations are combined into one vector (*delay vector*)

Starting from each location  $x_n$ , combine  $x_n$  with m subsequent locations if they were observed at a time interval v

$$\boldsymbol{x}_{n} = (x_{n}^{1}, x_{n}^{2}) \text{ location data points, index } n \in \{1, \dots, N\}$$
$$\boldsymbol{\delta}_{n} = [x_{n-(m-1)}^{1}, x_{n-(m-1)}^{2}, x_{n-(m-2)}^{1}, x_{n-(m-2)}^{2}, \dots, x_{n}^{1}, x_{n}^{2}]$$

For example:  $m = 2: \delta_n = [x_{n-1}^1, x_{n-1}^2, x_n^1, x_n^2]$ 



Note: Prior sampling with delay v is omitted for simplification

## Delay Embedding – Benefits

- Euclidean distance is a measure for similarity between subsequences
- Similar subsequences are close in embedding space
- Density is a measure for likelihood of a subsequence
- Mobility patterns can be extracted in terms of density



Learn mobility patterns from large amount of history data:

- Delay embedding to map mobility patterns to density
- Density estimation based on embedding space
  P(Xt = x, Xt-1, ..., Xt-(m-1))
- Derive conditional distribution

 $P(X_{t} = x | X_{t-1}, ..., X_{t-(m-1)}) = \alpha P(X_{t} = x, X_{t-1}, ..., X_{t-(m-1)})$ 

Predict location given the last m – 1 locations (current context):

 Maximization of probability density to obtain most likely location (MLL problem)

```
x^* = \underset{X}{\operatorname{argmax}} P(X_t = x, X_{t-1}, \dots, X_{t-(m-1)})
```

What about m=2?

```
Assuming (m-1)-th order Markov process •
```

RSITÄT ZU LÜBECK

MATIONSSYSTEMI

#### **Density Estimation**



Kernel Density Estimation

Optimization problem:

Minimize distance between estimated and unknown underlying distribution (AMISE, asymptotic mean integrated square error)

UNIVERSITÄT ZU LÜBECK

#### **Gaussian Mixture Models**



 $P(\mathbf{x}) = \sum \omega_m N(\mathbf{x} | \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$  $m \in M$ 



### **Online Kernel Density Estimation**

- Incremental can be updated as new data arrives
- Uses compression to keep memory footprint small

Christoph Heinz, Kernel Density Estimation over Data Streams, Dissertation Philipps-Universität Marburg, **2007** 

Matej Kristan, Aleš Leonardis, and Danijel Skočaj. 2011. Multivariate online kernel density estimation with Gaussian kernels. Pattern Recogn. 44, 10-11, 2630-2642, **2011** 



Master Thesis Jonas Lüthke, TUHH, 2013

IM FOCUS DAS LEBEN 27

# Solving MLL: Mode Finding



- Use hill-climbing search to find position of maximum
- Starting points?



Miguel Á. Carreira-Perpiñán. 2000. Mode-Finding for Mixtures of Gaussian Distributions. IEEE Trans. Pattern Anal. Mach. Intell. 22, 11, 1318-1323, **2000** 

IM FOCUS DAS LEBEN 28

Master Thesis Jonas Lüthke, TUHH, 2013

### Starting Points for Maxima Search

- Define search region around last observed location
- If radius large enough, all relevant maxima are found





#### **Summary - Prediction**

• Delay embedding:

Map mobility patterns to density

• Density estimation:

Assigns probability to each possible location sequence

• Mode finding:

Searches the most likely future location



#### **Test Results**

Varied *m*, fixed v = 6min:



Accurate predictions are more uniformly distributed for m = 3 and m = 5.



#### **Test Results**

Varied v, fixed m = 3:



Accurate predictions are increasingly clustered as v increases.



IM FOCUS DAS LEBEN 32

## **Test Result Analysis**

- Algorithm is based on sequential correlation in data (delay embedding)
- Locations in taxi data only correlated if part of same trip
- For each trip the client defines new destination
- Recurring similar location sequences only observed when limiting time span to average trip time
- Else prediction falls back to m = 2

#### Similar Approaches:

- Song et al. Markov predictor
- Scellato et al. Nonlinear predictor

L. Song, D. Kotz, R. Jain, and X. He, Evaluating location predictors with extensive with mobility data, In Proc. IEEE Computer and Communications Societies, pp. 1414-1424, **2004** 

S. Scellato, M. Musolesi, C. Mascolo, V. Latora, and A. T. Campbell, NextPlace: a spatio- temporal prediction framework for pervasive systems, In: Proc. Pervasive Computing, **2011** 

#### Duality Transformation: Avoid 3D-Rtrees?

- A linear trajectory in two-dimensional space can be transformed into a point in another *dual* two-dimensional space
- Trajectory:  $x(t) = vt + a \rightarrow Point: (v,a)$
- Embedding in more dimensions
- All queries will need to be transformed into the dual space



#### **Time Parameterized Queries**



10 a 8 d 6 b e 4 С the query q at time 1 2 x axis 8 10 2 0 4 6

v axis

- At time 1 b would be the nearest neighbor, after that time the results expire and d would be the new nearest neighbor
- Time Parameterized Query



#### Time Parameterized queries (TP queries)

- Whenever a query is issued, a TP returns:
  - Actual result that satisfies the corresponding spatial query.
  - Validity period/expiration time of the result.
  - Change that cause the expiration of the results
- Can be used for prediction


## **Time-Parameterized Data Structures**

- The Time-Parameterized R-tree (TPR-tree) consists of:
  - Minimum bounding rectangles (MBR)
  - Velocity bounding rectangles (VBR)
- A bounding rectangle with MBR & VBR is guaranteed to contain all its moving objects as long as they maintain their velocity vector



- Optimization: Minimize area of the bounding rectangle
- Time-Parameterized Bounding Rectangles (TPBRs) for answering TP queries



### Indexing Past, Present, and Future

- A unified index structure for both past, present, and future data
- Makes use of the partial-persistent R-tree for past data and the TPR-tree for current and future data



Katerina Raptopoulou, Michael Vassilakopoulos, and Yannis Manolopoulos.. Efficient processing of past-future spatiotemporal queries. In Proc. ACM Symposium on Applied Computing (SAC '06). ACM, pp. 68-72, **2006** 

IM FOCUS DAS LEBEN 38

Location-aware Query Processing and Optimization: A Tutorial, Mohamed F. Mokbel, Walid G. Aref

# Outline

- Location-aware Environments
- Location-aware Snapshot Query Processing
- Location-aware *Continuous* Query Processing
- Scalable Execution of Continuous Queries
- Location-aware Query Optimizer
- Uncertainty in Location-aware Query Processing



#### Approaches

- Straightforward Approach
  - Abstract the continuous queries to a series of snapshot queries evaluated periodically (and possibly incrementally)
- Result Validation
- Result Caching
- Result Prediction
- Incremental Evaluation



## **Result Validation**

- Associate a *validation* condition with each query answer
- Valid time (t):
  - The query answer is valid for the next t time units
- Valid region (R)
  - The query answer is valid as long as you are within a region *R*



- It is challenging to maintain the computation of valid time/region for querying *moving objects*
- Once the associated validation condition expires, the query will be *reevaluated*



## Caching the Result

- *Observation:* Consecutive evaluations of a continuous query yield very similar results
- Idea: Upon evaluation of a continuous query, retrieve more data that can be used later
- K-NN query
  - Initially, retrieve more than k
- Range query

VERSITÄT ZU LÜBECK

- Evaluate the query with a larger range
- How much do we need to pre-compute?
- How do we do re-caching?





Location-aware Query Processing and Optimization: A Tutorial, Mohamed F. Mokbel, Walid G. Aref

## Predicting the Result

- Given a future trajectory movement, the query answer can be pre-computed in advance
- The trajectory movement is divided into N intervals, each with its own query answers A<sub>i</sub>



Nearest-Neighbor Query

- The query is evaluated once (as a snapshot query). Yet, the answer is valid for longer time periods
- Once the trajectory changes, the query will be reevaluated



## **Incremental Evaluation**

- The query is evaluated only once. Then, only the *updates* of the query answer are evaluated
- There are two types of updates.
  *Positive* and *Negative* updates
- Only the objects that cross the query boundary are taken into account
- Need to continuously listen for notifications that someone crosses the query boundary









# Outline

- Location-aware Environments
- Location-aware Snapshot Query Processing
- Location-aware Continuous Query Processing
- Scalable Execution of Continuous Queries
  - Location-aware Centralized Database Systems
  - Location-aware Distributed Database Systems
  - Location-aware Data Stream Management Systems
- Location-aware Query Optimizer
- Uncertainty in Location-aware Query Processing



## Queries as Data – Motivation



Location-aware Query Processing and Optimization: A Tutorial, Mohamed F. Mokbel, Walid G. Aref

Continuous queries last for long times at the server side

- → While a query is active in the server, other queries will be submitted
- □ Shared execution among multiple queries

Should we index data OR queries?

- Data and queries may be stationary or moving
- → Data and queries are of large size
- → Data and queries arrive to the system with very high rates
- **Treat data and queries similarly**

Queries are coming to data OR data are coming to queries?

- → Both data and queries are subjected to each other
- Join data with queries



## Main Concepts (Cont.)



Evaluating a large number of concurrent continuous spatiotemporal queries is abstracted as a spatio-temporal join between moving objects and moving queries

UNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME

#### Location-aware Data Stream Management Systems

- Only *significant* objects are stored in-memory
- An object is considered significant if it is either in the query area or the cache area



- Due to the query and object movements, a stored object may become *insignificant* at any time
- Larger cache area indicates more storage overhead and more accurate answer



- The first k objects are considered an initial answer
- K-NN query is reduced to a circular range query

However, the query area may shrink or grow



*K* = 3





Location-aware Query Processing and Optimization: A Tutorial, Mohamed F. Mokbel, Walid G. Aref

#### • Query Load Shedding

- Reduce the cache area
- Possibly reduce the query area
- Immediately drop insignificant tuples
- Intuitive and simple to implement

#### Object Load Shedding

- Objects that satisfy less than k queries are *insignificant*
- Lazily drop insignificant tuples
- *Challenge I:* How to choose *k*?
- Challenge II: How to provide a lower bound for the query accuracy?





*K* = 2



# **Tutorial Outline**

- Location-aware Environments
- Jocation-aware Snapshot Query Processing
- J Location-aware Continuous Query Processing
- Scalable Execution of Continuous Queries
- Location-aware Query Optimization
- Uncertainty in Location-aware Query Processing



## Location-aware Query Optimization

- Spatio-temporal pipelinable query operators
  - Range queries
  - Nearest-neighbor queries
- Selectivity estimation for spatio-temporal queries/operators
  - Spatio-temporal histograms
  - Sampling
- Adaptive query optimization for continuous queries



## Spatio-temporal Query Operators

Existing Approaches are Built on Top of DBMS (at the Application Level)



55

## Spatio-temporal Query Operators





# Spatio-temporal Selectivity Estimation

- Estimating the selectivity of spatio-temporal operators is crucial in determining the best plan for spatio-temporal queries
- SELECT ObjectID FROM MovingObjects M WHERE Type = Truck INSIDE Region R









# Spatio-temporal Histograms

 Moving objects in D-dimensional space are mapped to 2Ddimensional histogram buckets





IM FOCUS DAS LEBEN 58

Location-aware Query Processing and Optimization: A Tutorial, Mohamed F. Mokbel, Walid G. Aref

## Spatio-temporal Histograms with Query Feedback

• Estimating the selectivity of spatio-temporal operators is crucial in determining the best plan for spatio-temporal queries





# Adaptive Query Optimization

- Continuous queries last for long time (hours, days, weeks)
  - Environment variables are likely to change
  - The initial decision for building a query plan may not be valid after a while
- Need continuous optimization and ability to change the query plan:
  - Training period: Spatio-temporal histogram, periodicity mining
  - Online detection of changes





# **Uncertainty in Moving Objects**

- Location information from moving objects is inherently inaccurate
- Sources of uncertainty:
  - Sampling. A moving object sends its location information once every t time units. Within any two consecutive locations, we have no clue about the object's exact location
  - Reading accuracy. Location-aware devices do not provide the exact location
  - Object movement and network delay. By the time that a certain reading is received by the server, the moving object has already changed its location



# **Uncertainty in Moving Objects**

• Historical data (Trajectories)

 $T_{a} + \epsilon_{a}$ 



Current data





# Error in Query Answer

• Range Queries



#### Nearest Neighbor Queries





## Representing Uncertain Data using Ellipses

- Given :
  - Start point
  - End point
  - Maximum possible speed  $\rightarrow$  Maximum traveling distance S
- If S is greater than the distance between the two end points, then the moving object may have deviated from the given route





#### Representing Uncertain Data using Cylinders

- Given:
  - Start and end points
- Constraint:
  - An object would report its location only if it is deviated by a certain distance r from the predicted trajectory





### Representing Uncertain Data in Road Networks

- Given:
  - Start and end points
- Constraints :
  - Deviation threshold r
  - Speed threshold v





## Querying Uncertain Data Uncertain Keywords

- KEYWORDS:
  - Probability: *possibly, definitely*
  - Temporal: sometimes, always
  - Spatial: *somewhere, everywhere*
- Examples:
  - What are the objects that are possibly sometimes within area R at time interval T?
  - What are the objects that definitely passed through a certain region?
  - Retrieve all the objects that are always inside a certain region
  - Retrieve all the objects that are sometimes definitely inside region R



## Querying Uncertain Data Uncertain Keywords (Cont.)



- Object O is definitely always in Q<sub>1</sub>
- Object O is possibly always in Q<sub>2</sub>

NIVERSITÄT ZU LÜBECK

- Object O is definitely sometimes in Q<sub>3</sub>
  - Object O is possibly sometimes in Q<sub>4</sub>

## Querying Uncertain Data Probabilistic Queries

- With each query answer, associate a probability that this answer is true
- The answer set of a query Q is represented as a set of tuples <ID, p> where ID is the tuple identifier and p is the probability that the object ID belongs to the answer set of Q
- Assumptions:
  - Objects can lie anywhere uniformly within their uncertainty region



## Querying Uncertain Data Probabilistic Range Queries



- Query Answer:
  - (B, 50%)
  - (C, 90%)
  - D
  - E
- (F, 30%)

## Querying Uncertain Data Probabilistic NN Queries

A



- Query Answer (k=1):
  - (C, *p*<sub>1</sub>)
  - (D, p<sub>2</sub>)
  - (E, p<sub>3</sub>)

UNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME

## Typicality Potential Fields (TyPoFs)



'Spieler vor dem Strafraum'


#### **Typicality Potential Fields**





J.R.J. Schirra: Bildbeschreibung als Verbindung von visuellem und sprachlichem Raum – Eine interdisziplinäre Untersuchung von Bildvorstellungen in einem Hörermodell. Dissertation. Infix, St. Augustin, **1994** 

## **Recap: Skyline Queries**

- Numeric space  $D = (D_1, ..., D_n)$ , larger values more preferable
- Two points, u dominates v (u > v), if  $- \forall D_i (1 \leq i \leq n), u.D_i \geq v.D_i$  $- \exists D_i (1 \leq j \leq n), u.D_i > v.D_i$
- Given a set of points S,

IVERSITÄT ZU LÜBECK

Skyline =  $\{u \mid u \in S \text{ and } u \text{ is not} \}$ dominated by any other point}

Example: C > B, C > D skyline = {A, C, E}





FOCUS DAS LEBEN

## Skylines on Uncertain Data

- Limitations of conventional methods
  - Aggregates may be misled by outliers
  - Data distribution is not captured
- Probabilistic skylines
  - Objects vs. instances
  - An instance has a probability to represent the object
  - An object has a probability to be in the skyline





Jian Pei, Bin Jiang, Xuemin Lin, and Yidong Yuan. 2007. Probabilistic skylines on uncertain data. In Proc. VLDB '07, 15–26, **2007**.

IM FOCUS DAS LEBEN

## A Probabilistic Skyline Model

- A set of objects S = {A, B, C}, instances of each with probability 0.5 to appear
- Probabilistic Dominance
  - Pr(A > C) = 3/4
  - Pr(B > C) = 1/2
  - $Pr((A > C) \lor (B > C)) = 1$



 $Pr(C \text{ is in the skyline}) \neq (1 - Pr(A > C)) \times (1 - Pr(B > C))$ 

#### Probabilistic dominance $\implies$ Probabilistic skyline



IM FOCUS DAS LEBEN

## **Skyline Probabilities**

• Possible world: W = <a<sub>i</sub>, b<sub>j</sub>, c<sub>k</sub>> (i, j, k = 1 or 2)

 $- Pr(W) = 0.5 \times 0.5 \times 0.5 = 0.125, \sum_{W \in \Omega} Pr(W) = 1$ 

- SKY(<a<sub>1</sub>, b<sub>1</sub>, c<sub>1</sub>>) = {a<sub>1</sub>, b<sub>1</sub>}
  - Objects A and B are in SKY(<a<sub>1</sub>, b<sub>1</sub>, c<sub>1</sub>>)
- B is in the skyline of possible worlds <a<sub>1</sub>, b<sub>1</sub>, c<sub>1</sub>>,
  <a<sub>1</sub>, b<sub>1</sub>, c<sub>2</sub>>, <a<sub>1</sub>, b<sub>2</sub>, c<sub>1</sub>>, and
  <a<sub>1</sub>, b<sub>2</sub>, c<sub>2</sub>>
  - $-Pr(B) = 4 \times 0.125 = 0.5$
- Pr(A) = 1, Pr(C) = 0





IM FOCUS DAS LEBEN

#### Problem Statement

- Skyline probability:  $Pr(U) = \sum_{U \in SKY(W)} Pr(W)$
- For object:  $Pr(U) = \frac{1}{|U|} \sum_{v \in U} \prod_{V \neq U} (1 \frac{|\{v \in V \mid v \succ u\}|}{|V|})$
- For instance:  $Pr(u) = \prod_{V \neq U} (1 \frac{|\{v \in V \mid v \succ u\}|}{|V|})$
- $Pr(U) = \frac{1}{|U|} \sum_{u \in U} Pr(u)$  Try to reduce V candidates
- p-skyline = {U |  $Pr(U) \ge p$ } for a given threshold p



IM FOCUS DAS LEBEN

## **Probabilistic Skyline Computation**

- Iteration: Bounding-Pruning-Refining
- Bounding
  - Bound Pr(u): lower bound  $Pr^{-}(u)$  and upper bound  $Pr^{+}(u)$

• Bound 
$$Pr(U)$$
:  $Pr(U) = \frac{1}{|U|} \sum_{u \in U} Pr(u)$ 

- Pruning
  - In *p*-skyline if lower bound  $Pr^{-}(U) ≥ p$
  - Not in *p*-skyline if upper bound  $Pr^+(U) < p$

## • Refining

- o Bottom-up method
- o Top-down method



IM FOCUS DAS LEBEN

## The Bottom-Up Method

- Sort instances of an object according to dominance relation such that their skyline probabilities are in descending order
- Two instances u<sub>1</sub> and u<sub>2</sub> ∈ U, if u<sub>1</sub> > u<sub>2</sub> then Pr(u<sub>1</sub>) ≥ Pr(u<sub>2</sub>)





IM FOCUS DAS LEBEN

## The Layer Structure

- layer-1: skyline of all instances
- layer-k (k > 1): skyline of instances except those at layer-1, ..., layer-(k-1)
- $\forall$  u at layer-k :  $\exists$  u' at layer-(k-1) : u'  $\succ$ u and Pr(u')  $\geq$  Pr(u)
- $max{Pr(u) | u is at layer-(k-1)} \ge max{Pr(u) | u is at layer-k}$
- Bounding example
  - $-\max{\Pr(u1), \Pr(u2)} \ge \max{\Pr(u3), \Pr(u4)} \ge \Pr(u5)$



IM FOCUS DAS LEBEN



## The Top-Down Method



 Build a partition tree for each object to organize partitions



IM FOCUS DAS LEBEN

#### **Partition Tree**



- Growing one level of the tree in each iteration
  - Choose one dimension in a round-robin fashion
  - Each leaf node is partitioned into two children nodes, each of which has half of instances
- Bound Pr(N<sub>max</sub>) and Pr(N<sub>min</sub>) of a partition N



IM FOCUS DAS LEBEN

# Summary

- Location-aware Environments
- Location-aware Snapshot Query Processing
- Location-aware *Continuous* Query Processing
- Scalable Execution of Continuous Queries
- Location-aware Query Optimizer
- Uncertainty in Location-aware Query Processing



IM FOCUS DAS LEBEN 84