# New Trends on Exploratory Methods for Data Analytics

***Davide Mottin, Matteo Lissandrini***,

*Yannis Velegrakis, Themis Palpanas*

db Trento

# Who we are

**Davide Mottin**

Graph Mining, Novel Query Paradigms, Interactive Methods

https://hpi.de/en/mueller/team/davide-mottin.html

**Matteo Lissandrini**

Knowledge Graphs , Novel Query Paradigms, Graph Mining

https://disi.unitn.it/~lissandrini

**Yannis Velegrakis**

Big Data Management & Analytics, Information Integration

https://velgias.github.io

**Themis Palpanas**

Data Series Indexing & Mining, Data Management, Data Analytics

http://www.mi.parisdescartes.fr/~themisp/

**Slides.** http://j.mp/DataExplore

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Big data – Easy value?

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Exploring



Traditional

On data

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Data exploration

Cleaning and profiling

Visualization

Analysis

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Data exploration software



Tableau: analysis and statistics

Trifacta:

OpenRefine: data preparation and cleanup

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Traditional data exploration methods

[Idreos et al., 2015]

Efficiently extracting knowledge from data
even if we do not know exactly what we are looking for

**SELECT** avg(system-stars)
**FROM** Universe
**WHERE** system-stars > 10
**GROUP BY** galaxy

Not easy for novices

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Declarative Exploratory methods

**SELECT** galaxy_name
**FROM** Universe.Galaxy

**SELECT** g.galaxy_name, SUM(s.stars) as st_s
**FROM** Universe.Galaxy  **AS** g
**JOIN** Universe.Systems **AS** s
**ON** g.galaxy_name = s.galaxy_name
**WHERE**
    g.st_s > 100B
    AND diameter > 100k AND diameter > 180k
    AND has_black_hole = TRUE
**GROUP BY** g.galaxy_name

Simple query (exploratory)

Complex query
(for data experts)

Over generic
100 billions results

Specific
Few results

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI  db Trento

# Examples as Exploratory Methods

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Historical perspective: Query-by-example

[Zloof et al. 1975]

Specify a query by example tables, or skeletons.

| Name | Stars | Diameter | Black_hole | Color | Life |
|------|-------|----------|------------|-------|------|
| P._  | > 10B | >100k    | TRUE       |       |      |
|      |       | <180k    |            |       |      |

- Intuitive GUI for simple queries
- SQL not required

- Restricted to SQL semantics
- Not example-based

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

db Trento

# Tutorial's goals

- Exploratory methods using examples
- Algorithms for retrieving data without using query languages
- Interactive methods and user-in-the-loop feedback
- Machine learning for adaptive, online methods

## But NOT

- Declarative query methods
- User interfaces and visualization
- Optimizations for fast data access
- Dynamic data

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Tutorial structure

Relational databases (25 min)

Textual data (10 min)

Graph and networks (25 min)

Machine learning (10 min)

Challenges and Remarks

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Example-based methods

- Query suggestion using examples
- Reverse engineering queries

- Entity extraction by example text
- Web table completion using examples
- Search by example

- Community-based Node-retrieval
- Entity Search
- Path and SPARQL queries
- Graph structures as Examples

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Where we are

Relational databases

Textual data

Graphs and networks

Challenges and Remarks

Machine learning

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis
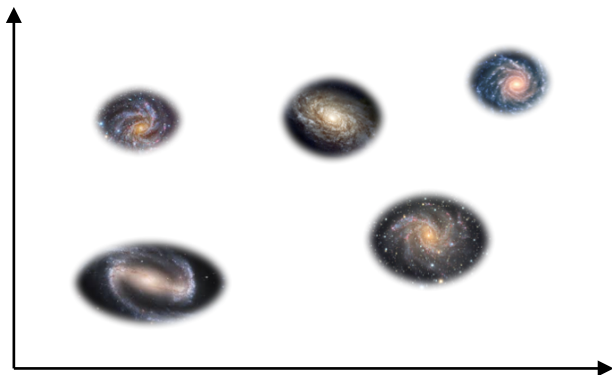
HPI   db Trento

# Reverse engineering queries (REQ)

Given a set of examples, find the query that generated that set of tuples

Example tuples



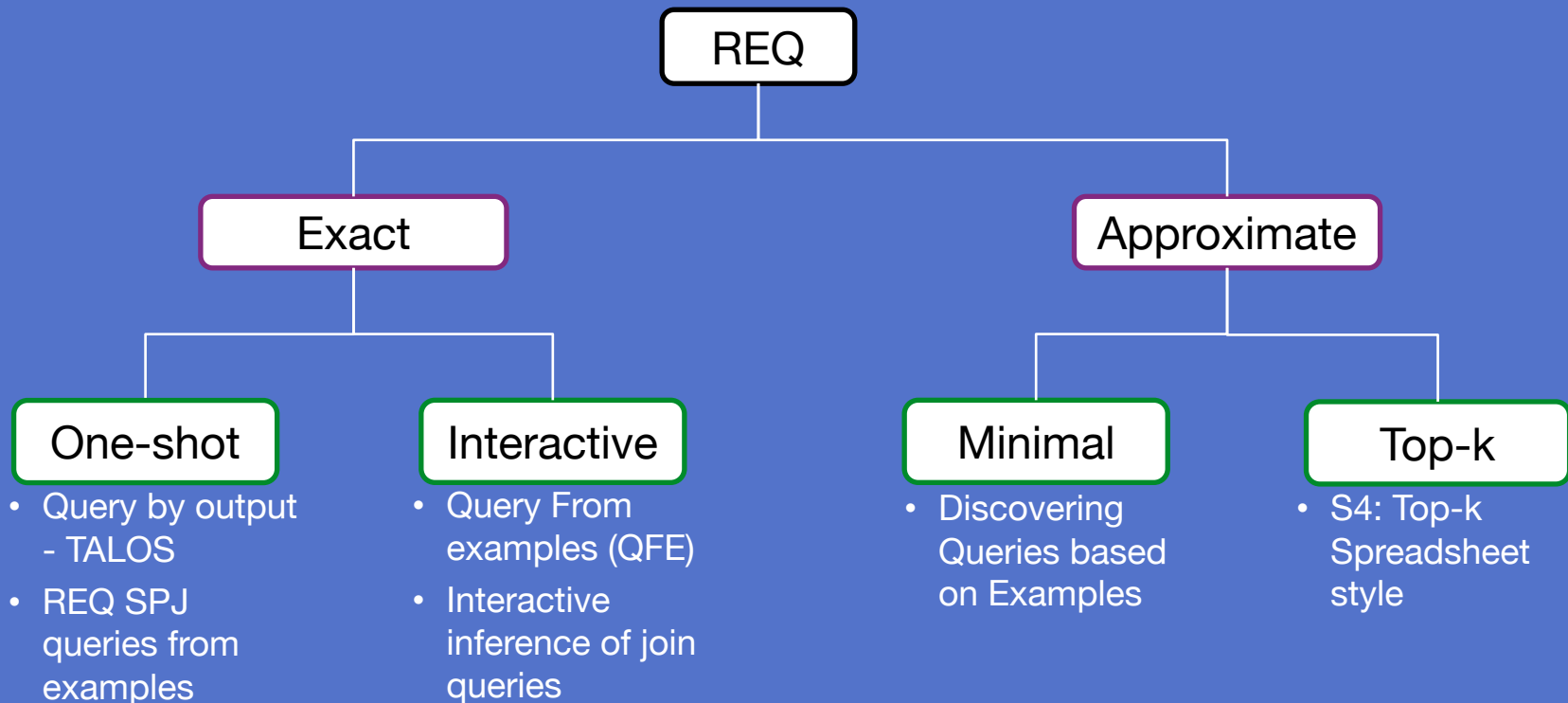How do you find such queries?

**SELECT** g.galaxy_name, SUM(s.stars) **AS** st_s
**FROM** Universe.Galaxy  **AS** g
**JOIN** Universe.System **AS** s
**ON** g.galaxy_name = s.galaxy_name
**WHERE**
    g.st_s > 100B
    AND diameter > 100k AND diameter > 180k
    AND has_black_hole = TRUE
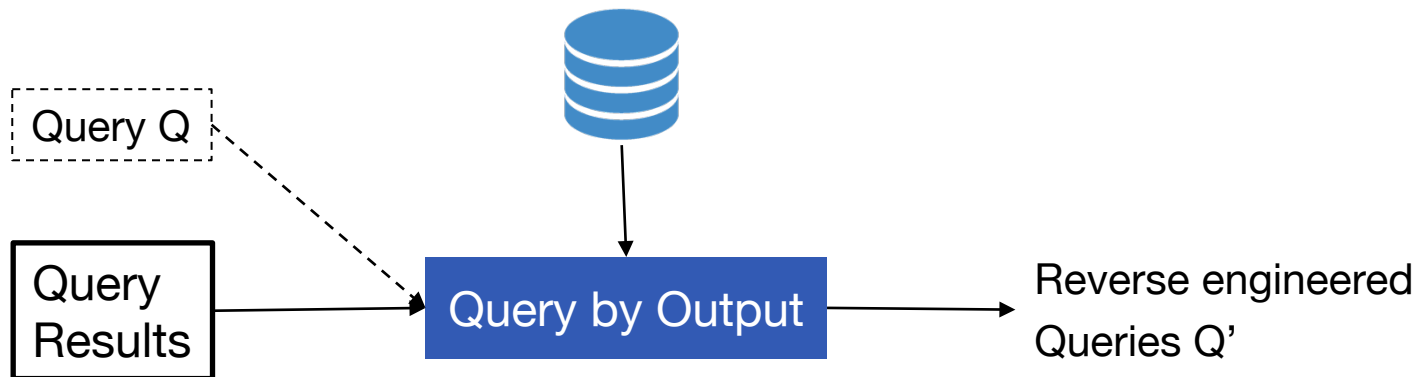**GROUP BY** g.galaxy_name

SELECT galaxy_name
FROM Universe.Galaxy

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Reverse engineering queries (REQ)



REQ

Exact

Approximate

One-shot
- Query by output - TALOS
- REQ SPJ queries from examples

Interactive
- Query From examples (QFE)
- Interactive inference of join queries

Minimal
- Discovering Queries based on Examples

Top-k
- S4: Top-k Spreadsheet style

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Query by Output - TALOS

**Main idea**: Find the set of queries that exactly return a set of examples

Query Q

Query Results → Query by Output → Reverse engineered Queries Q'

Two queries Q and Q' are instance equivalent on a database D, if the results of Q are the same of the results of Q'

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI dbTrento

# TALOS

| pID | name | country | weight | bats | throws |
|-----|------|---------|--------|------|--------|
| P1 | A | USA | 85 | L | R |
| P2 | B | USA | 72 | R | R |
| P3 | C | USA | 80 | R | L |
| P4 | D | Germany | 72 | L | R |
| P5 | E | Japan | 72 | R | R |

(a) Master

| pID | year | stint | team | HR |
|-----|------|-------|------|-----|
| P1 | 2001 | 2 | PIT | 40 |
| P1 | 2003 | 2 | ML1 | 50 |
| P2 | 2001 | 1 | PIT | 73 |
| P2 | 2002 | 1 | PIT | 40 |
| P3 | 2004 | 2 | CHA | 35 |
| P4 | 2001 | 3 | PIT | 30 |
| P5 | 2004 | 3 | CHA | 60 |

(b) Batting

| team | year | rank |
|------|------|------|
| PIT | 2001 | 7 |
| PIT | 2002 | 4 |
| CHA | 2004 | 3 |

(c) Team

Input

| B | PIT |
|---|-----|
| E | CHA |

Join table

$J = \texttt{Master} \bowtie \texttt{Batting} \bowtie \textit{Team}$

| | name | bat | throw | stint | HR | team |
|------|------|-----|-------|-------|-----|------|
| $t_1$ | A | L | R | 2 | 40 | PIT |
| $t_2$ | A | L | R | 2 | 50 | MT1 |
| $t_3$ | C | R | L | 2 | 35 | CHA |
| $t_4$ | D | L | R | 3 | 30 | PIT |
| $t_5$ | B | R | R | 1 | 73 | PIT |
| $t_6$ | B | R | R | 1 | 40 | PIT |
| $t_7$ | E | R | R | 3 | 60 | CHA |

Master ← Batting → Team

Join graph computation

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI  dbTrento

# TALOS

| | name | bat | throw | stint | HR | team | |
|---|---|---|---|---|---|---|---|
| $t_1$ | A | L | R | 2 | 40 | PIT | ✗ |
| $t_2$ | A | L | R | 2 | 50 | MT1 | ✗ |
| $t_3$ | C | R | L | 2 | 35 | CHA | ✗ |
| $t_4$ | D | L | R | 3 | 30 | PIT | ✗ |
| $t_5$ | B | R | R | 1 | 73 | PIT | ✓ |
| $t_6$ | B | R | R | 1 | 40 | PIT | ✓ |
| $t_7$ | E | R | R | 3 | 60 | CHA | ✓ |

**Idea**: treat the problem as a binary classification

1. **Strict**: all tuples must be captured
2. **At-Least-one**: one tuple for example must be captured



Decision tree

$$Gini(S) = 1 - (f_+^2 + f_-^2)$$

Positive and negative tuples in S

$$Gini(S_1, S_2) = \frac{(|S_1|Gini(S_1) + |S_2|Gini(S_2))}{|S_1| + |S_2|}$$

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# How complex is exact REQ?

Database $D$

Relational Operators:
$\sigma$ selection $\{=, \neq, \geq, \leq\}$
$\pi$ projection
$\bowtie$ natural join

$E^+$ Positive examples
$E^-$ Negative examples

→ REQ →

$Q$ such that results contain
- All positive examples
- No negative example

How difficult is to find:
A bounded size Q?  an unbounded Q?

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Complexity  - No parameters

| Operator | Unbounded Queries | Bounded Queries |
|---|---|---|
| $\pi$ | P | P |
| $\bowtie$ | P | NPC |
| $\sigma$ | P | NPC |
| $\sigma, \bowtie$ | P | NPC |
| $\pi, \sigma$ | NPC | NPC |
| $\sigma, \bowtie$ | DP | DP |
| $\pi, \sigma, \bowtie$ | DP | DP |

Only projections: Easy

Unbounded selections: Easy
Unbounded selections: HARD

Combination of operators:
HARD!!!

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Unbounded Select

| | A | B | C | D | E |
|---|---|---|---|---|---|
| ☑ | 1 | 2 | 3 | 4 | 5 |
| ☒ | 1 | 3 | 2 | 3 | 4 |
| | 2 | 4 | 4 | 1 | 3 |
| | 5 | 3 | 2 | 4 | 2 |
| ☒ | 4 | 2 | 3 | 1 | 2 |
| | 2 | 2 | 4 | 3 | 2 |
| ☑ | 1 | 1 | 2 | 1 | 5 |
| ☑ | 1 | 5 | 4 | 2 | 3 |

**Possible queries?**

$A = 1$    AND

$B \geq 1$    AND    $B \leq 5$    AND

$C \geq 2$    AND    $C \leq 4$    AND

$D \geq 1$    AND    $D \leq 4$    AND    $D \neq 4$

$E \geq 3$    AND    $E \leq 5$    AND    $E \neq 4$

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Bounded select

Reduction from Set Cover

NP-C

INPUT: Database D, Examples E, Query size k

OUTPUT: Does there exist a query satisfying D and E, of size at most k?

U = {1,2,3,4,5}     S = { {1,2,3}, {2,4}, {3,4}, {4,5} }

| | $S_1$ | $S_2$ | $S_3$ | $S_4$ |
|---|---|---|---|---|
| ☒ | 1 | 0 | 0 | 0 |
| ☒ | 1 | 1 | 0 | 0 |
| ☒ | 1 | 0 | 1 | 0 |
| ☒ | 0 | 1 | 1 | 1 |
| ☒ | 0 | 0 | 0 | 1 |
| ☑ | 0 | 0 | 0 | 0 |

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Complexity - Parameters

| | No param | Schema | Examples | No param | Query | Schema | Examples |
|---|---|---|---|---|---|---|---|
| $\pi$ | $P$ | - | - | $P$ | - | - | - |
| $\bowtie$ | $P$ | - | - | $NPC$ | $P$<br>$W[2]C$ | $P$<br>$FPT$ | $NPC$ |
| $\sigma$ | $P$ | - | - | $NPC$ | $P$<br>$W[2]C$ | $\{=\}:P,\{\neq\}:NPC$<br>$\{=\}:FPT$ | $P$<br>$FPT$ |
| $\sigma,\bowtie$ | $P$ | - | - | $NPC$ | $P$<br>$W[2]C$ | $\{=\}:P,\{\neq\}:NPC$<br>$\{=\}:FPT$ | $NPC$ |
| $\pi,\sigma$ | $NPC$ | $\{=\}:P,\{\neq\}:NPC$<br>$\{=\}:FPT$ | $P$<br>$W[3]C$ | $NPC$ | $P$<br>$W[3]C$ | $\{=\}:P,\{\neq\}:NPC$<br>$\{=\}:FPT$ | $NPC$ |
| $\pi,\bowtie$ | $DP$ | $P$<br>$W[1]H,$ co-$W[1]H$ | $DP$ | $DP$ | $P$<br>$W[2]H,$ co-$W[1]H$ | $P$<br>$W[1]H,$ co-$W[1]H$ | $DP$ |
| $\pi,\sigma,\bowtie$ | $DP$ | $\{=\}:P,\{\neq\}:NPC$<br>$\{=\}:W[1]H,$ co-$W[1]H$ | $DP$ | $DP$ | $P$<br>$W[3]H,$ co-$W[1]H$ | $\{=\}:P,\{\neq\}:NPC$<br>$\{=\}:W[1]H,$ co-$W[1]H$ | $DP$ |

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Interactive REQ – Query from Examples

[Li et al., 2015]

**Main idea**: Interactively remove candidate queries proposing a new set of query results from a modified database

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Database Refinement

[Li et al., 2015]

| eid | name | gender | dept | salary |
|-----|------|--------|------|--------|
| 1 | Alice | F | Sales | 3700 |
| 2 | Bob | M | IT | 4200 |
| 3 | Carol | F | Service | 3000 |
| 4 | Dave | M | IT | 5000 |

Results

| name |
|------|
| Bob |
| Dave |

Database
Refinement

| eid | name | gender | dept | salary |
|-----|------|--------|------|--------|
| 1 | Alice | F | Sales | 3700 |
| 2 | Bob | M | IT | 3000 ~~4200~~ |
| 3 | Carol | F | Service | 3000 |
| 4 | Dave | M | IT | 5000 |

✅

| name |
|------|
| Bob |
| Dave |

Result $R_1' = Q_1(D') = Q_3(D')$

| name |
|------|
| ~~Bob~~ |
| Dave |

Result $R_2' = Q_2(D')$

REQs =
- $Q_1 = \sigma_{gender=M}(D)$
- $Q_2 = \sigma_{salary>3700}(D)$
- $Q_3 = \sigma_{dept=IT}(D)$

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Cost model

Number of modified tables

Number of new result sets

$$cost(D') = \overline{edit(D, D') + \beta \cdot n} + \sum_{i=1}^{k} edit(R, R_i) + N \cdot \frac{edit(D, D')}{\mu} + \beta + \frac{2}{k} \sum_{i=1}^{k} edit(R, R_i)$$

DB cost | Results cost

Effort to examine D' | Effort to examine new results

**Current cost** | **Residual cost**

Main idea: Find a refined db D' and results $R_1, \dots R_k$ with:
1. Minimum number of results k
2. Minimum differences i the database
3. The query are balanced (less interactions)

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI  dbTrento

# Minimal Project Join REQ

**Main idea**: Find the set of queries that approximately return a set of examples



| | A | B | C |
|---|---|---|---|
| 1 | Mike | ThinkPad | Office |
| 2 | Mary | iPad | |
| 3 | Bob | | Dropbox |

- valid: every tuple is present in query results
- minimal: any removal in query tree gets to an invalid query

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Candidate Query Generation

- Use candidate network generation algorithm (Hristidis 2002)

|   | A | B | C |
|---|------|---------|---------|
| 1 | Mike | ThinkPad | Office |
| 2 | Mary | iPad | |
| 3 | Bob | | Dropbox |

**CQ₁**

Sales
A → Customer
B → Device
C → App

**CQ₂**

Owner
A → Employee
B → Device
C → App

**CQ₃**

Owner
A → Employee
B → Device
C → ESR

**CQ₄**

Owner
C → App
B → Device
ESR
A → Employee

**CQ₅**

Owner
A → Employee
B → Device
App
C → ESR

1. Generate join tree $J$
2. Generate mapping $\phi$
3. Check minimal:
   - Every leaf node contains a column that is mapped by an input column

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Validity verification

Naïve: check all candidate queries singularly if they return ALL examples

Better: exploit substructures in candidate queries for pruning

Best: adaptively select the substructures to have the min number of evaluations

NP-hard



$CQ_2$

Candidate query

Substructures

Sub 1

Sub 1 fails => $CQ_2$ invalid

Sub 2

Sub 1 fails => Sub 2 fails

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI  db Trento

# Minimal Project Join REQ

**Main idea: Allow missing rows/columns and rank the k best queries**

**Partial** query table

S4

Output: Top-k PJ Queries

Sales
Products    Customers
Name    First Name    Last Name

Sales
Products    Customers
Name    Last Name    City
Name

| | A | B | C |
|---|---|---|---|
| 1 | John | Smith | Xbox |
| 2 | Jill | Hans | Surface |

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI  dbTrento

# Ranking score

Linear combination of row score and column score

$$\frac{\alpha * score_{row}(Q) + (1 - \alpha) * score_{col}(Q)}{|Q|}$$

- $\alpha = 1$ penalizes missing rows
- $\alpha = 0$ penalizes missing columns

**Row score**



**Column score**

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# S4 Optimizations

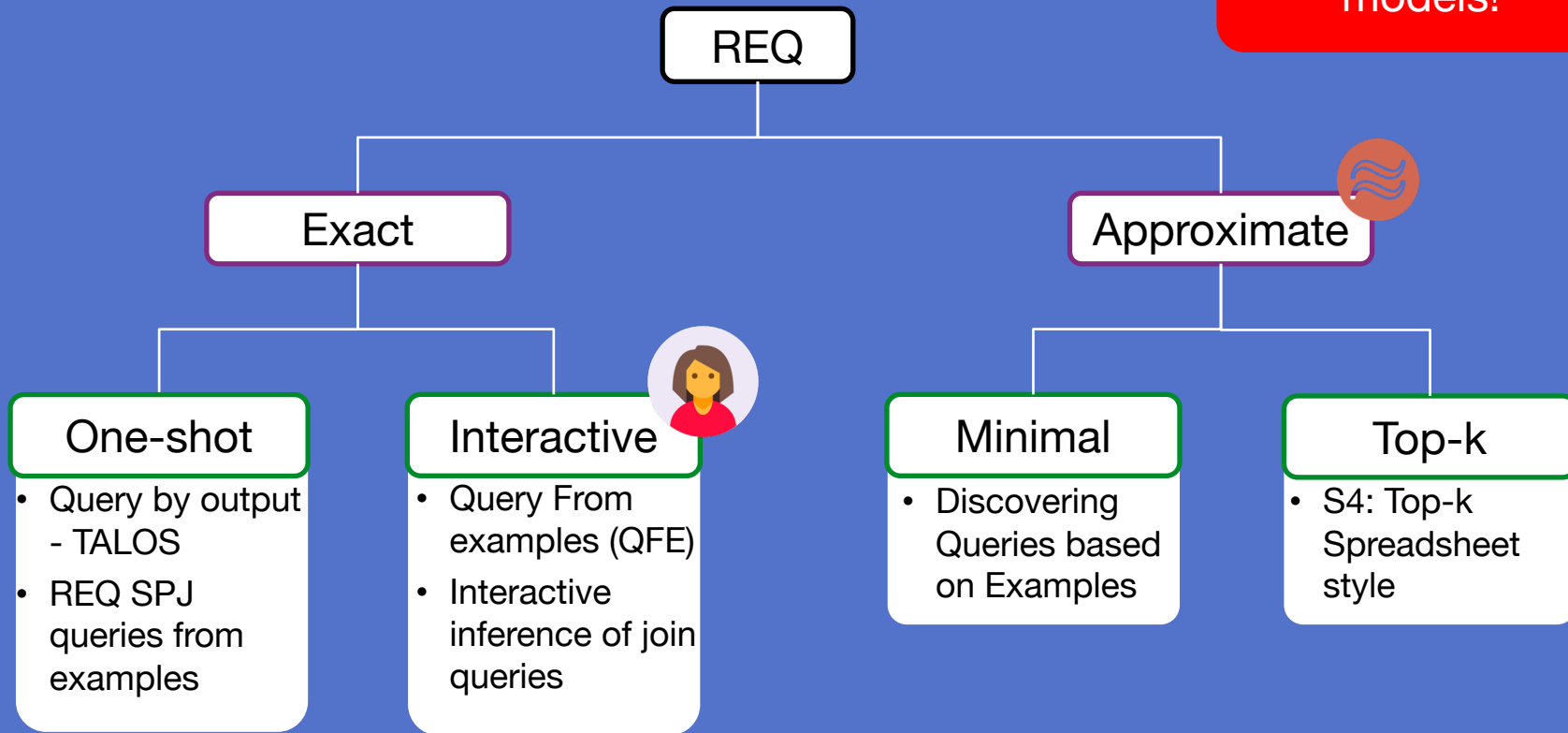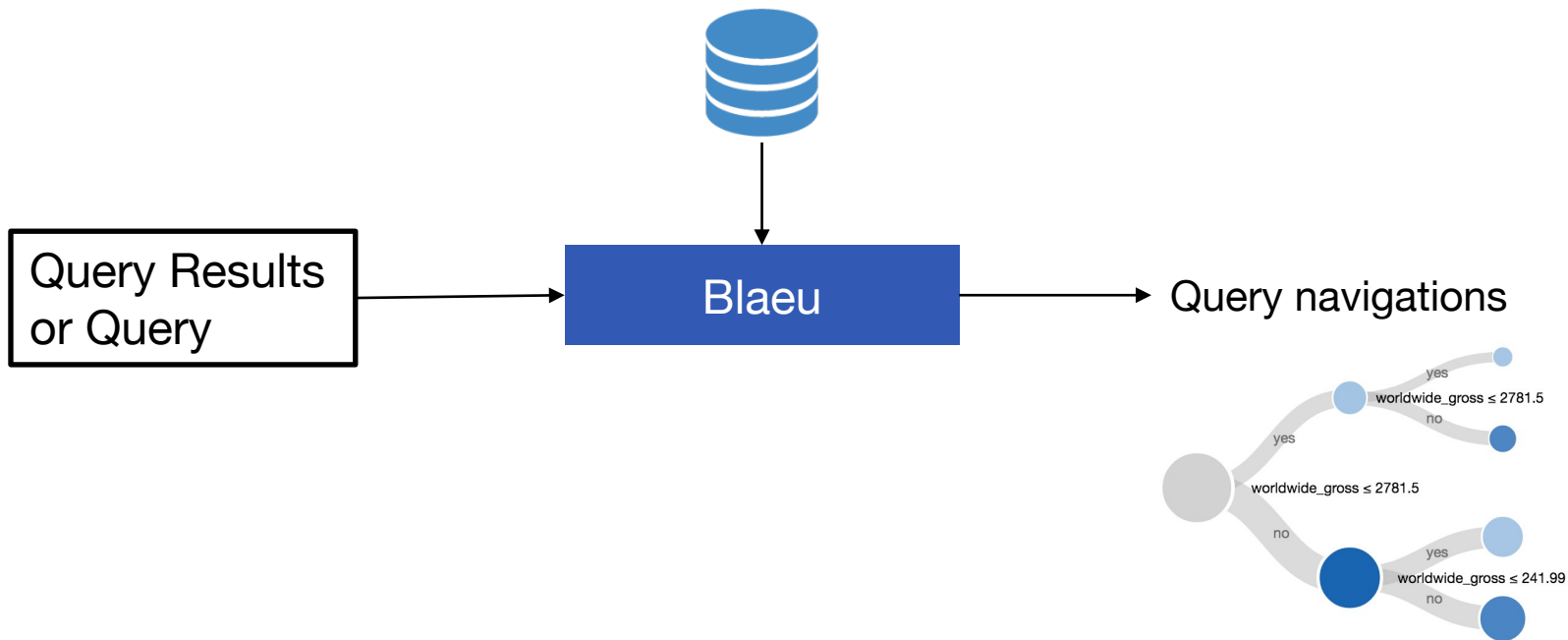| Upper bound | Row score is always bounded by the column score<br>(row containment is more restrictive)<br><span style="color:red">Exploit inverted indexes on columns/rows</span> |
| --- | --- |
| Early termination | Stop when current upper bound score is less than the k-th ranked evaluated query<br><span style="color:red">Scan queries on decreasing upper bound</span> |
| Caching | Reuse common subparts in the candidate queries |

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Reverse engineering queries (REQ)

**Lack of user models!**

REQ

**Exact**

**Approximate**

**One-shot**
- Query by output - TALOS
- REQ SPJ queries from examples

**Interactive**
- Query From examples (QFE)
- Interactive inference of join queries

**Minimal**
- Discovering Queries based on Examples

**Top-k**
- S4: Top-k Spreadsheet style

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Examples for query suggestion: Blaeu

[Sellam et al., 2016]

**Main idea**: Allow interactive navigation of the query space in a hierarchy

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Examples for query suggestion: Blaeu

[Sellam et al., 2016]

Query results



Attribute 2

Attribute 1

$$u: DB \rightarrow \{-1, 1\}, U(Q) = \sum_{t \in Q} u(t)$$

User utility

Given a result of an example query Q, explore the data through data maps = partitions

**Output**: Set of query refinements

**Problem**: User utility is unknown

- Cluster analysis for result exploration
- Zoom and projection operations
- User model

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI db Trento

# Examples for query suggestion: Blaeu

[Sellam et al., 2016]

$$u: DB \rightarrow \{-1, 1\}, U(C) = \sum_{t \in C} u(t)$$

Unknown User utility

Find the partition $\mathcal{C} = \{C_1, \ldots, C_n\}$ of the results of Q such that exists $C_j \in \mathcal{C}: U(C_j) > U(Q)$
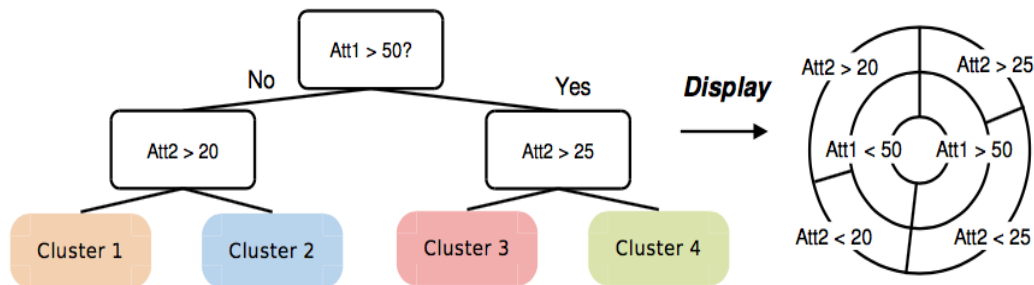
**Solution**: interesting tuples are close to each other within a maximum separation threshold $\theta(\mathcal{C})$
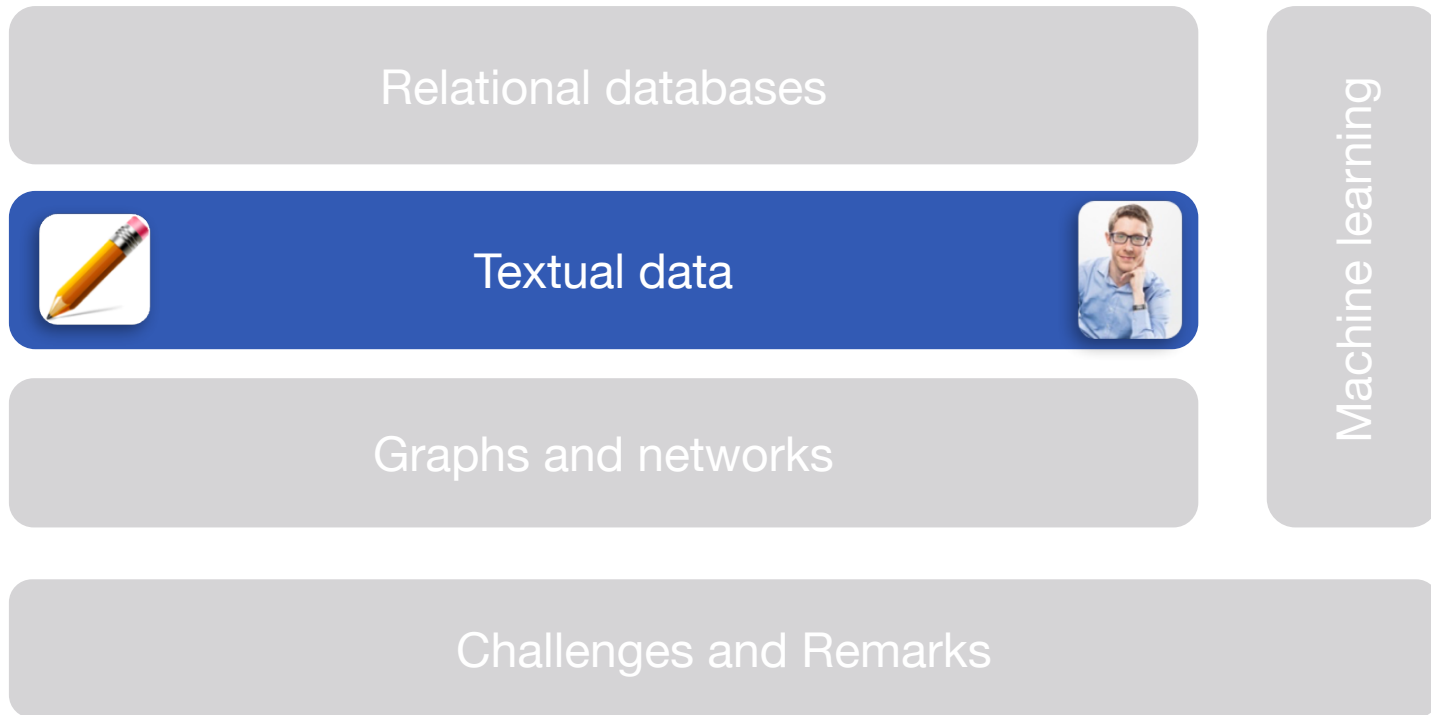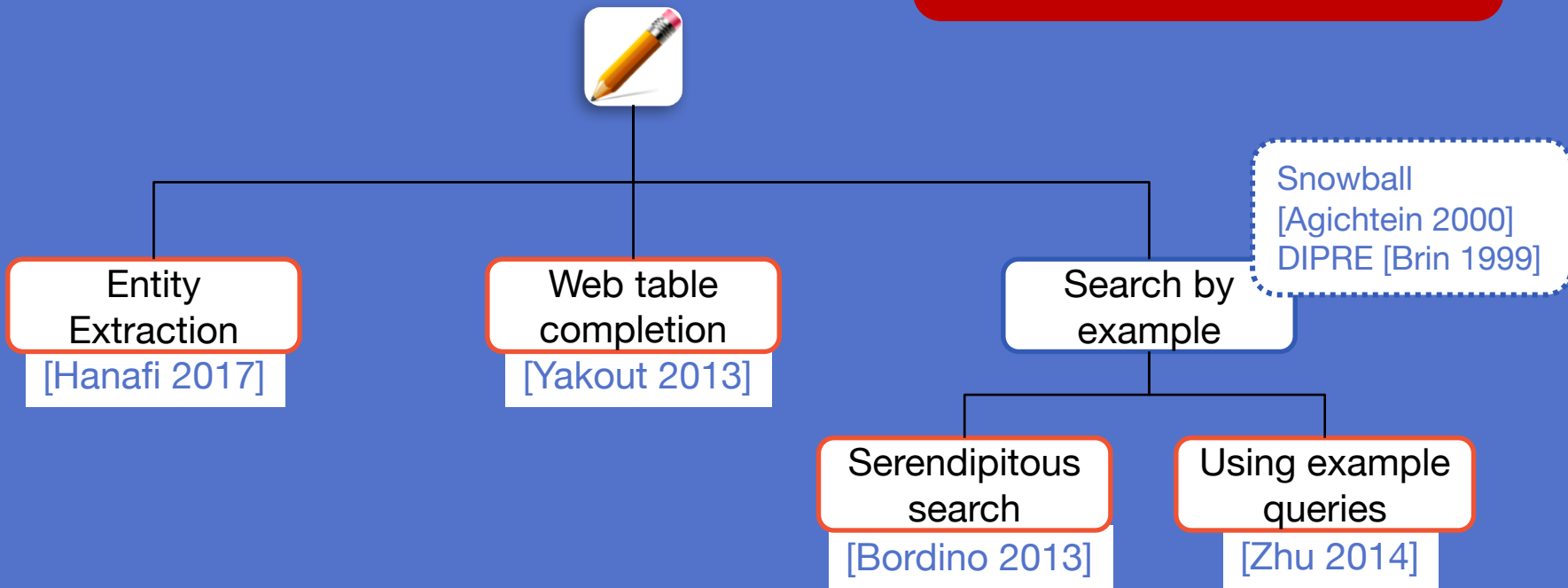
### Detect clusters (k-medoid)

### Organize clusters

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Where we are

Relational databases

Textual data

Graphs and networks

Challenges and Remarks

Machine learning

**VLDB 2017** tutorial    D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Examples for textual data

Few methods for textual data using examples



Entity Extraction
[Hanafi 2017]

Web table completion
[Yakout 2013]

Search by example

Snowball [Agichtein 2000]
DIPRE [Brin 1999]

Serendipitous search
[Bordino 2013]

Using example queries
[Zhu 2014]

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Entity extraction by-example (SEER)

[Hanafi et al., 2017]

**Main idea:** Create rules to extract wanted information from documents using examples

definition) increased 9.6 percent, the number of murders increased 6.2 percent, aggravated assaults increased 2.3 percent, the number of rapes (revised definition) rose 1.1 percent, and robbery violations were up 0.3 percent.
Violent crime increased in all but two city groupings. In cities with populations from 50,000 to 99,999 inhabitants, violent crime was down 0.3 percent, and in cities with 500,000 to 999,999 in population, violent crime decreased 0.1 percent. The largest increase in violent crime, 5.3 percent, was noted in cities with 250,000

SEER

**Output**: Extraction rules

P: Percentage = 1.0     = 1.0
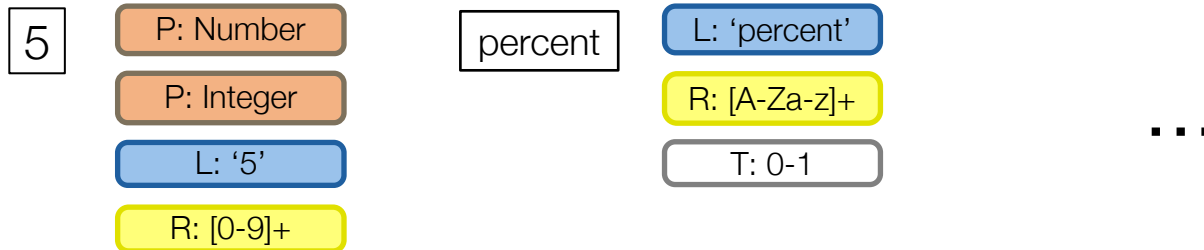
D: {5, 6} = 0.4   D: {percent, %} = 0.4   = 0.4

R: [0-9]+ = 0.2   D: {percent, %} = 0.4   = 0.3
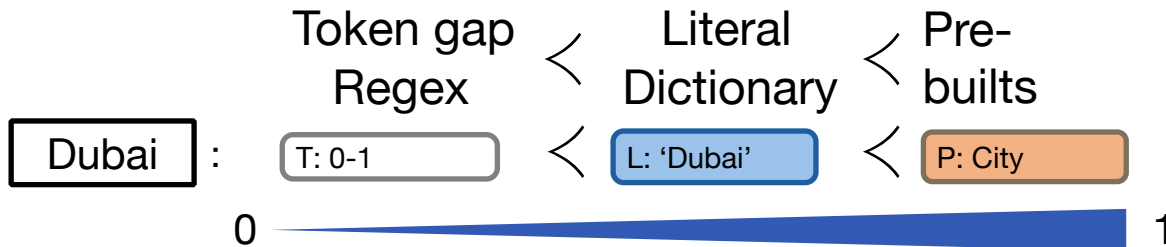
D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI   db Trento

# Learning rules

Example:    **5 percent up**

1.    **Enumerate** possible primitives per example token

| 5 |

P: Number
P: Integer
L: '5'
R: [0-9]+

| percent |

L: 'percent'
R: [A-Za-z]+
T: 0-1

...

2.    **Assign** scores to primitives

Token gap
Regex    <    Literal
Dictionary    <    Pre-
builts

| Dubai | :    T: 0-1    <    L: 'Dubai'    <    P: City

0                                                                     1

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis         HPI dbTrento

# Learning rules (cont'd)

[Hanafi et al., 2017]

## 3. Generate rules

Example:  **5 percent**

Tokens:     5          percent

Tree:

P: Percentage = 1.0

L: 'percent' = 0.4
R: [A-Za-z]+ = 0.2

L: '5' = 0.4

R: [0-9]+ = 0.2    L: 'percent' = 0.4
R: [A-Za-z]+ = 0.2

Rule:    R: [0-9]+ = 0.2    L: 'percent' = 0.4

Example:  **6%**

P: Percentage = 1.0          L: '%' = 0.4

L: '6' = 0.4                 R: symbols = 0.2

R: [0-9]+ = 0.2              L: '%' = 0.4

R: symbols = 0.2

## 4. Merge

Intersection:   [**5 percent, 6%**]

P: Percentage = 1.0

D: {5, 6} = 0.4 ——— D: {percent, %} = 0.4

R: [0-9]+ = 0.2 ——— D: {percent, %} = 0.4

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI  dbTrento
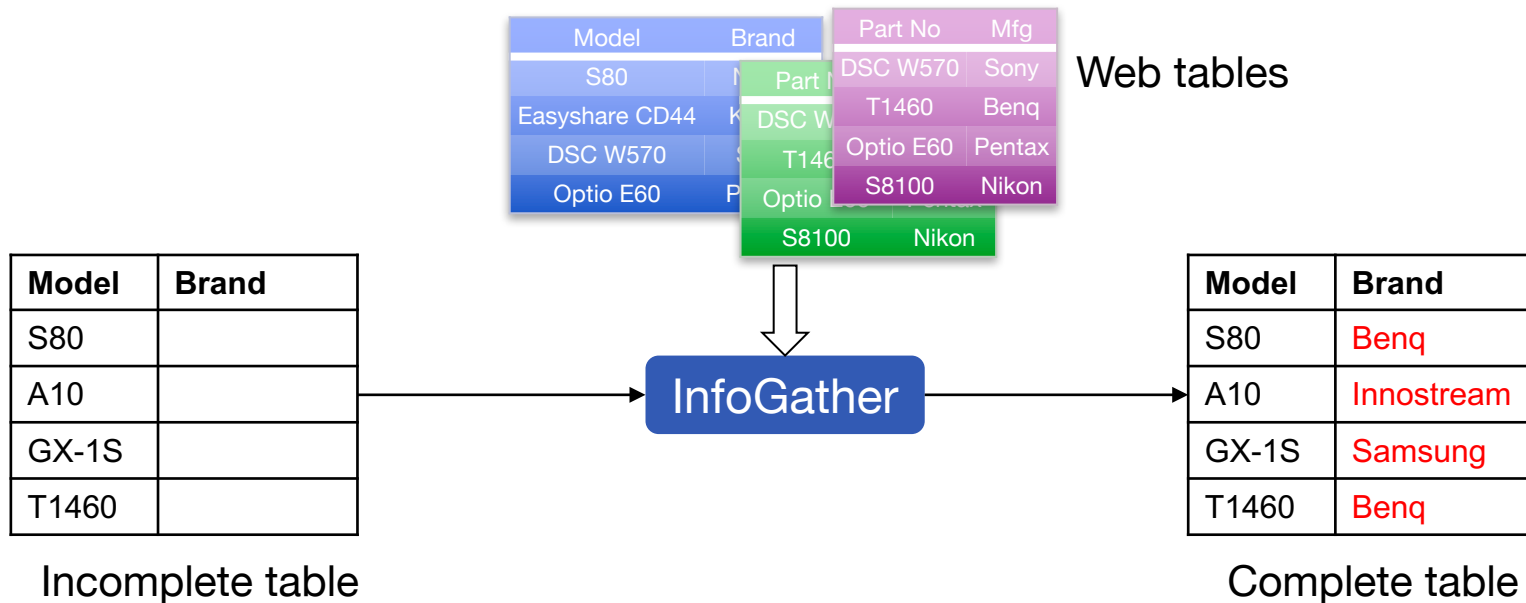
# Web tables completion (InfoGather)

[Yakout et al., 2012]

**Main idea:** Complete tables using partial information about tuples



Web tables
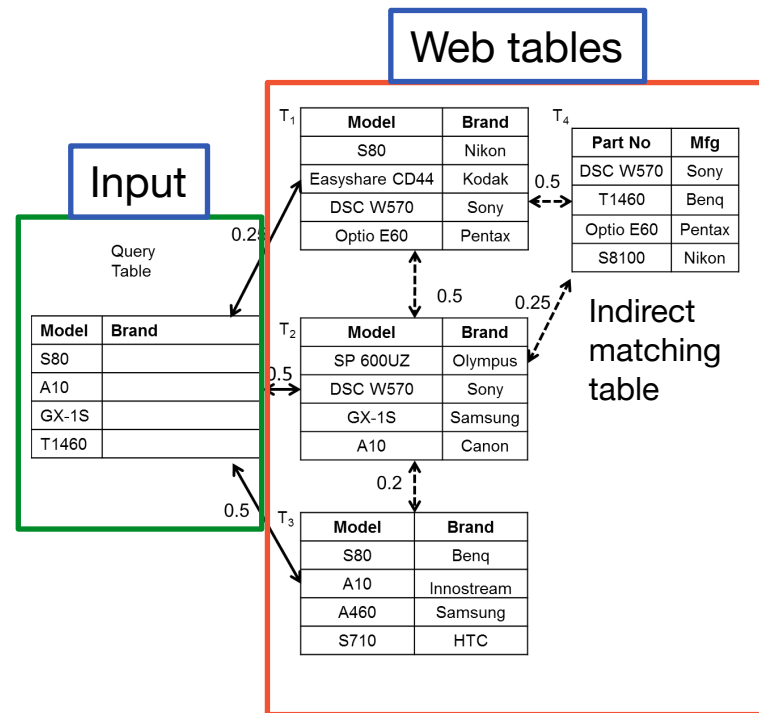
| Model | Brand |
|-------|-------|
| S80 | |
| A10 | |
| GX-1S | |
| T1460 | |

Incomplete table

InfoGather

| Model | Brand |
|-------|-------|
| S80 | Benq |
| A10 | Innostream |
| GX-1S | Samsung |
| T1460 | Benq |

Complete table

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

[Yakout et al., 2012]

## Direct Match Approach (DMA)

- Traditional schema matching techniques using the attribute names and the values in the column

$$S_{DMA}(T) = \begin{cases} \dfrac{|T \cap_K Q|}{\min(|Q|, |T|)} & if\ Q.A \approx T.B \\ 0 & otherwise \end{cases}$$

Web tables

Input

T₁

| Model | Brand |
|---|---|
| S80 | Nikon |
| Easyshare CD44 | Kodak |
| DSC W570 | Sony |
| Optio E60 | Pentax |

T₄

| Part No | Mfg |
|---|---|
| DSC W570 | Sony |
| T1460 | Benq |
| Optio E60 | Pentax |
| S8100 | Nikon |

0.5

0.25

Query Table

| Model | Brand |
|---|---|
| S80 | |
| A10 | |
| GX-1S | |
| T1460 | |

0.5

0.25

Indirect matching table

T₂

| Model | Brand |
|---|---|
| SP 600UZ | Olympus |
| DSC W570 | Sony |
| GX-1S | Samsung |
| A10 | Canon |

0.5

0.2

0.5

T₃

| Model | Brand |
|---|---|
| S80 | Benq |
| A10 | Innostream |
| A460 | Samsung |
| S710 | HTC |

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Ranking tables using PageRank

- **PageRank**

- **Personalized PageRank (PPR)**

$$\pi_u(v)$$
$$= \epsilon\, \delta_u(v) + (1 - \epsilon) \sum_{\{w\,|\,(w,v)\in E\}} \pi_u(w)\alpha_{w,v}$$

Adjacency matrix

- **Topic Sensitive Pagerank (TSP)**

$$\pi_{\vec{\beta}}(v) = \epsilon\, \vec{\beta} + (1 - \epsilon) \sum_{\{w\,|\,(w,v)\in E\}} \pi_{\vec{\beta}}(w)\alpha_{w,v}$$
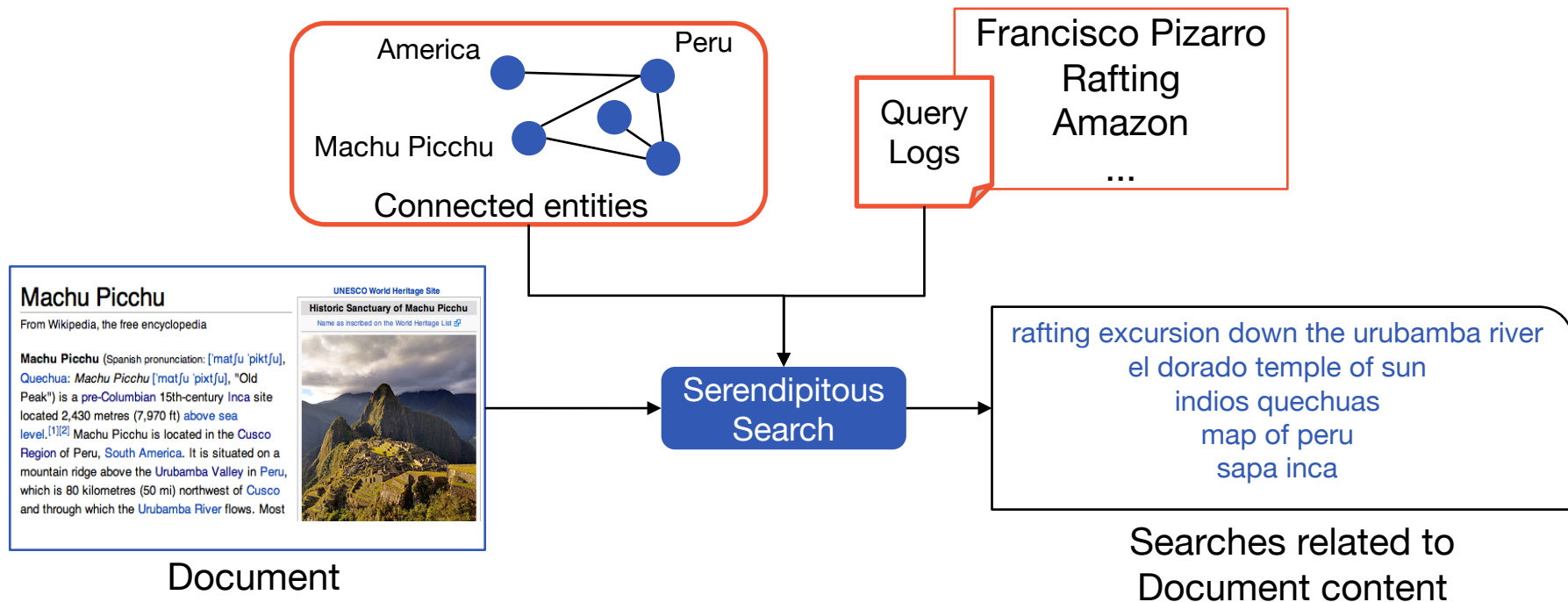
Topic vector

Query Table

Nodes ➔ Web Tables
Edges ➔ Tables Similarity

Topic weight ➔ DMA score

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI  db Trento

# Serendipitous search

[Bordino et al., 2013]

**Main idea:** Use related entities and query logs to find serendipitous searches

America — Peru

Machu Picchu

Connected entities

Query Logs

Francisco Pizarro
Rafting
Amazon
...

## Machu Picchu

From Wikipedia, the free encyclopedia

**Machu Picchu** (Spanish pronunciation: [ˈmatʃu ˈpiktʃu], Quechua: *Machu Picchu* [ˈmatʃu ˈpixtʃu], "Old Peak") is a pre-Columbian 15th-century Inca site located 2,430 metres (7,970 ft) above sea level.[1][2] Machu Picchu is located in the Cusco Region of Peru, South America. It is situated on a mountain ridge above the Urubamba Valley in Peru, which is 80 kilometres (50 mi) northwest of Cusco and through which the Urubamba River flows. Most

UNESCO World Heritage Site
**Historic Sanctuary of Machu Picchu**
Name as inscribed on the World Heritage List

**Document**

**Serendipitous Search**

rafting excursion down the urubamba river
el dorado temple of sun
indios quechuas
map of peru
sapa inca

Searches related to
Document content

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI    db Trento

# Find queries using entity-query graph

[Bordino et al., 2013]

Entity Nodes

Query Nodes



**Idea**: Run Personalized PageRank on entity-query graphs

**Query-flow graph with entity nodes**

Three types of arcs:

**1. query to query:**

$$w_Q(q_i \rightarrow q_j) = w_{QFG}(q_i \rightarrow q_j)$$

**2. entity to query**

Frequency-based approach

$$w_{EQ}(e \rightarrow q) = \frac{f(q)}{\sum_{q_i \mid e \in X_E(q_i)} f(q_i)}$$

**3. entity to entity**

The more queries entities share the higher the probability

$$w_E(e_u \rightarrow e_v) = 1 - \prod_{i=1,\ldots,r} \left(1 - \boxed{p_{q_{i_s} \rightarrow q_{i_t}}(e_u \rightarrow e_v)}\right)$$

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Search by multiple examples

**Main idea:** Document examples are used to find topics

Chuck Norris
Arnold Schwarzenegger

→ Search by examples →

Action Movies
- Mission impossible
- Die Hard
- …

Action Actors
- Bruce Willis
- Tom Cruise
- …

…

Related topics and documents

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI  db Trento

# Nearest neighbor approach

**Main Idea:**

The similarity is an aggregation over the distances between document $D_i$ and its nearest query example

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Where we are

Relational databases

Textual data

**Graphs and networks**

Challenges and Remarks

Machine learning

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI db Trento

# Graphs



**Fact Graph**

**Ontology Tree**

Arnold Schwarzenegger

is A → Person

actedIN

is A

subClassOf

Terminator

| Release | 1984 |
| Budget | $6.4M |
| Length | 1h 48m |

Actor

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Graphs



**RDF**

**(subject,predicate,object)**

```
(Arnold_Schwarzenegger,isA,Person)
      (Actor, subClassOf, Person)
(Arnold_Schwarzenegger, actedIn, Terminator)
```

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Exemplar Queries

**Input:** $Q_e$, an example ***element*** of interest

**Output:** set of elements in the desired result set

> Nodes/Entities
> Edges/Facts
> Structures

## Exemplar Query Evaluation

- **evaluate** $Q_e$ in a database D, finding a sample $S$
- **find** the set of elements $A$ similar to $S$ given a *similarity relation*

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Exemplar Queries

**Input:** $Q_e$, an example ***element*** of interest

**Output:** set of elements in the desired result set

> Nodes/Entities
> Edges/Facts
> Structures

### Exemplar Query Evaluation

- **evaluate** $Q_e$ in a database D, finding a sample $S$
- **find** the set of elements $A$ **similar** to $S$ given a *similarity relation*
- **[OPTIONAL]** return only the subset $A^R$ that are **relevant**

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# SIMILARITY

**Nodes**

**Structures**

### Connectivity

Mediator Nodes
[Ruchansky'15]

Clusters
[Perozzi'14]

### Properties

Entity Search
[Metzger'13,
Sobczak'15]

### Queries

Path Queries
[Bonifati'15]

SPARQL
[Arenas'16]

### (Edge-)Labels

Entity Tuples
[Jayaram'15]

Graph Structures
[Mottin'14]

**CHALLENGE: DISCOVER USER PREFERENCE**

**CHALLENGE: EFFICIENT SEARCH**

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# The Minimum Wiener Connector Problem

[Ruchansky, et al., 2015]

**Model:** **Unlabeled Undirected Graph**

**Query:** A set of **Nodes Q**

**Similarity:** Shortest-Path **distance**

**Output:** A Set of **Connector Nodes H**
"*explains*" connections in **Q**

Connectors:
Nodes with **HIGH** closeness
to **ALL** the inputs

Similar to a Steiner-Tree but
**overall pairwise distances** are optimized

Case: Infected Patients
→ Culprit/Other Infected

Case: Target Audience
→ Influencers

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI db Trento

# The Minimum Wiener Connector Problem

[Ruchansky et al., 2015]

**Model:** **Unlabeled Undirected Graph**

**Query:** A set of **Nodes Q**

**Similarity:** Shortest-Path **distance**

**Output:** A Set of **Connector Nodes H** **minimize** the sum of **pairwise shortest-path-distances** between nodes in the **connector H**

W=1+2+1 =4

W=1+1+1 = 3

Sometimes The Best Solution is NOT A Tree

NP-Hard

**Called: Wiener Index.**

*tradeoff between size*

*and average distance*

$$min \sum_{(u,v) \in H} d(u, v)$$

d(u, v) is the shortest-path distance

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI db Trento

# Approximate minimum Wiener Index Connector

CHOOSE $r$ & $\lambda \in \left[ 1, \ \log_{(1+\beta)} |V| \right]$

Approximated with
**Edge-Weighted SteinerTree**

All Pairwise Distances

     ⮑ **Distances from a root r**

**Enumerate** Candidate Solutions
for $r \in Q$ & $\lambda$
and **keep best**

Measure distance in H

     ⮑ **Precomputed distance in G**

Edge Weights

$$w(u, v) = \quad \lambda + \frac{max\{d_G(r,u), d_G(r,v)\}}{\lambda}$$

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Focused Clustering and Outlier Detection

[Perozzi et al., 2014]

**Model:** **Unlabeled Undirected Graph**
**with Node Attributes**

**Query:** A set of **Nodes Q**

**Similarity: Attribute Values & Connectivity**
(*to be inferred*)

**Output: Clusters** of Nodes: Dense & Coherent
**+Cluster Outliers**

Case: Target Users → Community with same interests

Case: Products→ Co-purchased products with similar features

**PhD NYC** English Google

**PhD NYC** Greek SAP

College Paris Dutch Google

**College NYC** English Google

**PhD NYC** Italian IBM

**PhD NYC** French SAP

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI dbTrento

# Focused Clustering and Outlier Detection

[Perozzi et al., 2014]

**TASK: Infer "FOCUS" , important attributes**

**attribute weights β**

PhD
NYC
English
Google

PhD
NYC
French
SAP

→

0.5
0.5
0
0

**1. Set of similar pairs, PS  (from Q)**

**2. Set of dissimilar pairs, PD (random sample)**

**3. Learn a distance metric between PS and PD**
*( Distance Metric Learning, inverse Mahalanobis distance: Xing, et al 2002)*

College
Paris
Dutch
Google

PhD
NYC
English
Google

PhD
NYC
Greek
SAP

**College**
**NYC**
English
Google

PhD
NYC
Italian
IBM

PhD
NYC
French
SAP

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

db Trento

# Focused Clustering and Outlier Detection

[Perozzi et al., 2014]

**TASK: Extract Clusters on Focused Graph**

**attribute weights β -> Edge Weight**

**LOCAL** clusters

**1. Find Starting Set of Candidates**

1.a Drop low-weight edges

1.b Extract **Strongly Connected Component** $C_1, C_2, ...$

**2. Grow Clusters around Candidates**

2.a Compute conductance of **C**: $\phi^{(w)}$ **(C, G)**

2.b Select node to add to **C'**: **best improvement to $\Delta\phi^{(w)}$ (C,C')** *(greedy)*

2.c Prune Underperforming nodes

**3. Detect Outliers:** High *unweighted* conductance

Seed

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# SIMILARITY

**Nodes**

**Structures**

**Connectivity**

Mediator Nodes
[Ruchansky'15]

Clusters
[Perozzi'14]

✓

**Properties**

Entity Search
[Metzger'13,
Sobczak'15]

**Queries**

Path Queries
[Bonifati'15]

SPARQL
[Arenas'16]

**(Edge-)Labels**

Entity Tuples
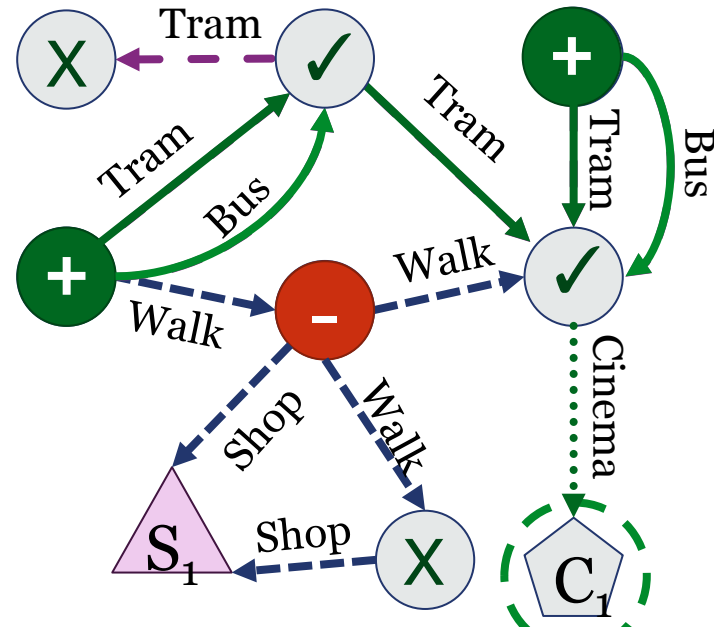[Jayaram'15]

Graph Structures
[Mottin'14]

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI   db Trento

# iQBEES: Entity Search by Example

[Metzger et al., 2013, Sobczak et al., 2015]

**Model:** Knowledge Graph

**Query:** A set of Entities Q

**Similarity:** shared semantic properties

**Output:** A Set of Similar Entities
*ranked*

Case: Products → Find Similar Products

Case: Social Media → User recommendation

Entity 1:

Entity 2:

?

?

?

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Maximal Aspects

?x type BodyBuilder

?x type AmericanActor

**Adding any aspect**
**→ E(A)={Arnold}**

?x type AmericanActor

?x type GovernorCalifornia

**Include**
**Typical Types**

**Prune generic aspects**

?x hasHeight 1.88m

?x type Entity

**use most specific type**

**Rank Set of aspects**

?x type AmericanActor

?x actedIn TheExpendables

?x type ActionActor

**REPEATABLE**
*Update Q*

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# SIMILARITY

**Nodes**

**Structures**

| Connectivity | Properties | Queries | (Edge-)Labels |
|---|---|---|---|
| Mediator Nodes [Ruchansky'15] | Entity Search [Metzger'13, Sobczak'15] | Path Queries [Bonifati'15] | Entity Tuples [Jayaram'15] |
| Clusters [Perozzi'14] ✓ | ✓ | SPARQL [Arenas'16] | Graph Structures [Mottin'14] |

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Learning Path Queries on Graphs

[Bonifati et al., 2015]

**Model:** **Edge Labeled Graph**

**Query:** **2 sets of Entities Q$^+$ , Q$^-$**
**Positive, Negative**

**Similarity:** **common path query (RegExp)**
(bus|tram)*Cinema

**Output:** **A Set of Nodes Satisfying**
**some paths(Q$^+$) but NOT paths(Q$^-$)**

Case: Proteins→ Similar interactions/co-expression

Case: Tasks Initiator→ Similar Processes/Behaviours



**MONADIC: only starting nodes**
*extensible to*
**BINARY/ N-ARY** : path from X to Y

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Learnability of Path Queries

**Query: $Q^+$ & $Q^-$ (Positive & Negative examples)**

**Consistecy:** $\forall v \in Q^+ . paths_G(v) \not\subseteq paths_G(Q^-)$

**Consistency Check: PSPACE-complete**

**Enumerate Paths Up to Fixed distance**

**For paths of Length N K = 2 X N +1**

## 1. Selecting the Smallest Consistent Paths

Infinite Paths? **Fix maximal length $K$ but...**

When to use **Kleene star * ?**

$$C \,|\, (A \cdot B \cdot C) \rightarrow (A \cdot B)^* \cdot C$$

## 2. Generalize SCP

a. Construct Prefix-Tree Acceptor

b. Generalize into DFA with Merge



**PTA**

**DFA**

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Reverse engineering SPARQL queries

[Arenas et al., 2016]

**Model:** **Knowledge Graph**

**Query:** **Set of *ANSWERS*** *

**Similarity:** **common** AND/OPT/FILTER **query**

**Output:** **A SPARQL QUERY/RESULT**

Case: Open Data → Query Unknown Schema

Case: Novice User → Avoid SPARQL



| | ?e1 | ?e2 |
|---|---|---|
| **M1** | Mexico | Spanish |
| **M2** | Haiti | |
| **M3** | Jamaica | English |

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Reverse engineering SPARQL queries

[Arenas et al., 2016]

**Query:**     **Set of _Variable Mappings_**

|      | $?X$  | $?Y$            | $?Z$         |
|------|-------|-----------------|--------------|
| **M1** | John  |                 |              |
| **M2** | Mary  | mary@email.eu   |              |
| **M3** | Lucy  |                 | Roses Street |

$(?X, \texttt{type}, \texttt{Person})$ $\quad ?X \neq \texttt{me}$

    OPT        OPT

$(?X, \texttt{email}, ?Y)$       $(?X, \texttt{addr}, ?Z)$

Enumerate all possible SPARQL queries satisfied by the mappings

INTRACTABLE
$\Sigma_2^p-\text{complete}$
$\mathbf{coNP}-\text{complete}$

Build tree-shaped SPARQL queries IMPLIED by the mappings

    D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI   db Trento

# Reverse engineering SPARQL queries

[Arenas et al., 2016]

**Query:** **Set of _Variable Mappings_ Ω**

$\Omega$

|    | $?X$ | $?Y$ | $?Z$ | $?W$ |
|----|------|------|------|------|
| **M1** | a1 |  |  |  |
| **M2** | a2 | b2 |  |  |
| **M3** | a3 |  | c3 |  |
| **M4** | a4 | b4 | c4 | d4 |

{M1,M2,M3,M4} $?X$

{M2,M4} $?Y$    {M3,M4} $?Z$

{M4} $?W$

$D$

| | |
|---|---|
| **M1** | $(\mathtt{a1},\mathtt{t},\mathtt{P})$ |
| **M2** | $(\mathtt{a2},\mathtt{t},\mathtt{P})(\mathtt{a2},\mathtt{e},\mathtt{b2})$ |
| **M3** | $(\mathtt{a3},\mathtt{t},\mathtt{P})(\mathtt{a3},\mathtt{a},\mathtt{c3})$ |
| **M4** | $(\mathtt{a4},\mathtt{t},\mathtt{P})(\mathtt{a4},\mathtt{e},\mathtt{b4})$ $(\mathtt{a4},\mathtt{a},\mathtt{c4})$ $(\mathtt{b4},\mathtt{d},\mathtt{d4})$ |

Greedy: keep just enough to cover all variables

$?X$    $?X$

$?Y$    $?Z$ $?Y$    $?Z$

$?W$    $?W$

$(?X, \mathtt{t}, \mathtt{P})\ ?X$

OPT        OPT

$(?X, \mathtt{e}, ?Y)\ ?Y$    $(?X, \mathtt{a}, ?Z)\ ?Z$

OPT

$(?Y, \mathtt{d}, ?W)(\mathtt{b4}, \mathtt{d}, ?W)\ (?Y, \mathtt{d}, \mathtt{d4})\ ?W$

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI    db Trento

# SIMILARITY

**Nodes**

**Structures**

## Connectivity

Mediator Nodes
[Ruchansky'15]

Clusters
[Perozzi'14]

✓

## Properties

Entity Search
[Metzger'13,
Sobczak'15]

✓

## Queries

Path Queries
[Bonifati'15]

SPARQL
[Arenas'16]

✓

## (Edge-)Labels

Entity Tuples
[Jayaram'15]

Graph Structures
[Mottin'14]

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Exemplar Queries

**Model:** **Knowledge Graph**

**Input:** **Example Structure**

**Similarity:** **Isomorphism/Simulation**

**Output: A set of Graphs**

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Computing exemplar queries [Mottin et al., 2014]

NP-complete
(subgraph isomorphism)

$O(|V|^4)$ (simulation)



Q    Sample    A1

Labels at distance 1

A2

Pruning technique:
- Compute the neighbor labels of each node

$W_{n,a,i} = \{n_1 | l(n_1, n_2) = a \lor \in N_{i-1}(n)\}$

- Prune nodes not matching query nodes neighborhood labels

- Apply iteratively on the query nodes

v neighborhood = {(B,1)}

$\not\subseteq$

No Match

u neighborhood = {(A,1)}

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Computing exemplar queries [Mottin et al., 2014]

NP-complete
(subgraph isomorphism)

$O(|V|^4)$ (simulation)



Sample    A1

A2

Approximation:
- Nodes closed to the sample are more important
- Use Personalized PageRank with a weighted matrix

$$\boldsymbol{v} = (1 - c)A\boldsymbol{v} + c\boldsymbol{p}$$

- Weight edges: <u>frequency of the edge-label</u>

$$I(e_{ij}^{\ell}) = I(\ell) = \log \frac{1}{P(\ell)} = -\log P(\ell)$$

$$P(\ell) = \frac{|E^{\ell}|}{|E|}$$

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI   dbTrento

# Ranking results



[Mottin et al., 2014]

$$\rho(n_s, n) = \lambda \mathcal{S}(n_s, n) + (1 - \lambda)\boldsymbol{v}[n]$$

**Combination of two factors**
1. Structural: similarity of two nodes in terms of neighbor relationships
2. Distance-based: the PageRank already computed

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Graph query by example (GQBE)

**Model:** **Knowledge Graph**

**Input:** **Entity Tuples**

**Similarity: Isomorphism**

**Output: A set of Tuples**

In GQBE Input is a set of (disconnected) entity mention tuples

Q = (Google, S. Mateo)

Results =
(Yahoo, S. Clara)
(CBS, New York)

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# GQBE: Maximum Query Graph

[Jayaram et al., 2015]

$Q = (v_1, v_2)$



Maximum Query Graph

Answer graph

1. Find the maximum query graph
   • Graph with <u>M edges</u> having the <u>maximum weight</u>

2. Answers subgraph-isomorphic to the query graph    **NP-hard**

3. Return top-k

**Answer score:**
• Sum of query graph weights
• Similarity match between edges in the answer and the query (shared nodes take extra credit)

$$\text{match}(e, e') = \begin{cases} \frac{w(e)}{|E(u)|} & \text{if } u = f(u) \\ \frac{w(e)}{|E(v)|} & \text{if } v = f(v) \\ \frac{w(e)}{min(|E(u)|, |E(v)|)} & \text{if } u = f(u), v = f(v) \\ 0 & \text{otherwise} \end{cases}$$

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Multiple query tuples

**Subgraphs of Maximum Query graph**

v₁ — v₂

v₁ — v₂    v₁ — v₂    v₁ — v₂

v₁ — v₂

**Maximum Query Graph is Very Large**

Preserve the query connectivity

**Find answers using a lattice obtained removing edges from the union graph**

**GQBE finds answers for multiple query tuples**

1. **Compute a re-weighted union graph of the individual query graphs**

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# SIMILARITY

**Nodes**

**Structures**

| Connectivity | Properties | Queries | Structures |
|---|---|---|---|
| Mediator Nodes [Ruchansky'15]  Clusters [Perozzi'14] | Entity Search [Metzger'13, Sobczak'15] | Path Queries [Bonifati'15]  SPARQL [Arenas'16] | Entity Tuples [Jayaram'15]  Graph Structures [Mottin'14] |

**Do not Include User Feedback**

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Where we are

Relational databases

Textual data

Graphs and networks

Challenges and Remarks

Machine learning

**VLDB 2017** tutorial

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Online exploration of datasets

**Main idea:** Learn the items to show online as more points are acquired

Two ways of learning: passive and active



Passive

Active

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# MindReader

**Main idea:** learn an implicit query from user examples and optional scores

Searching "mildly overweighted" patients

• The doctor selects examples by browsing patient database

• The examples have **"oblique"** correlation

• We can "guess" the implied query



✓ : good

✓✓ : very good

Weight

Height

q

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Learning an ellipsoid distance

**Euclidean**



**weighted Euclidean**



**generalized ellipsoid distance**



**Weighted distance matrix**

$$D(x, q) = (x - q)^\top M(x - q)$$

**Implicit query**

$$D(x, q) = \sum_{j}^{n} \sum_{k}^{n} m_{jk}(x_j - q_j)(x_k - q_k)$$

Learn the query minimizing the penalty = weighted sum of distances between query point and sample vectors

$$minimize \sum_{i} (x_i - q)^\top M(x_i - q)$$

$$subject\ to \quad \det(M) = 1$$

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Learning the distance

■ Query point is moved towards "good" examples — Rocchio formula in IR



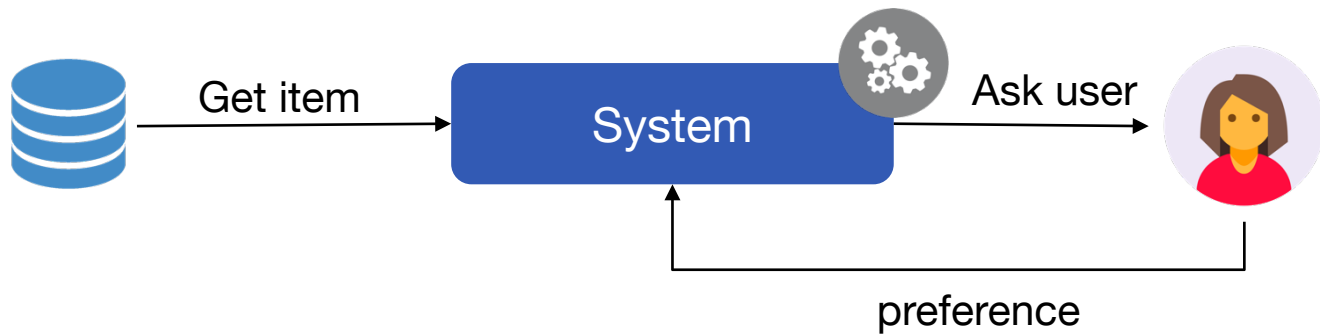$Q_0$: query point

● : retrieved data

✓ : relevance judgments

$Q_1$: new query point

### Learning can be done online!!!

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Active learning for online query systems

**Main idea:** the system "query" the user to understand her preferences



Get item → System → Ask user

preference

Learn unknown preferences and minimize the number of questions to the user

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Learning unknown preferences

[Vanchinathan et al., 2015]

**Problem**: Find a set S that maximize the user preference within a budget (e.g., number of interactions)

S (intended user set)

User preferences

$$\arg\max \sum_{v \in S} pref(v)$$

$$\text{subject to } Cost(S) \leq budget$$

Cost for the set S

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Background: Gaussian processes

[Bishop et al., 2006]

**Idea**: Model the user preferences as a Gaussian Process

A Gaussian Process (GP) is an infinite set of variables, any subset of this is Gaussian
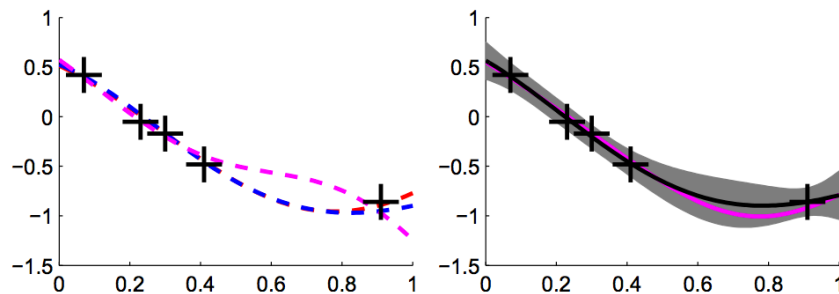
$$P(\mathbf{f}|\Sigma, \mu) = |2\pi\Sigma|^{\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{f} - \mu)^{\top}\Sigma^{-1}(\mathbf{f} - \mu)\right)$$

Gaussian prior

Specified only by mean and covariance

Given observations $\{x, y\}_{i=1}^{n}$ over an unknown function f drawn from a Gaussian prior, the posterior is Gaussian

$$P(\mathbf{f}|\mathbf{y}) \propto \int d\mathbf{x}\, P(\mathbf{f}, \mathbf{x}, \mathbf{y})$$

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# GP-Select

**Algorithm 1** GP-SELECT

**Input:** Ground Set $\mathbf{V}$, kernel $\kappa$ and budget $B$
Initialize selection set $S$
**for** $t = 1, 2, \ldots, B$ **do**
    **Model Update:**
        $[\mu_{t-1}(\cdot), \sigma_{t-1}^2(\cdot)] \leftarrow$ GP-Inference$(\kappa, (S, y_{\{1:t-1\}}))$
    **Item Selection:**
        Set $v_t \leftarrow \underset{v \in \mathbf{V}/\{v_{1:t-1}\}}{\mathrm{argmax}} \mu_{t-1}(v) + \beta_t^{1/2} \sigma_{t-1}(v)$
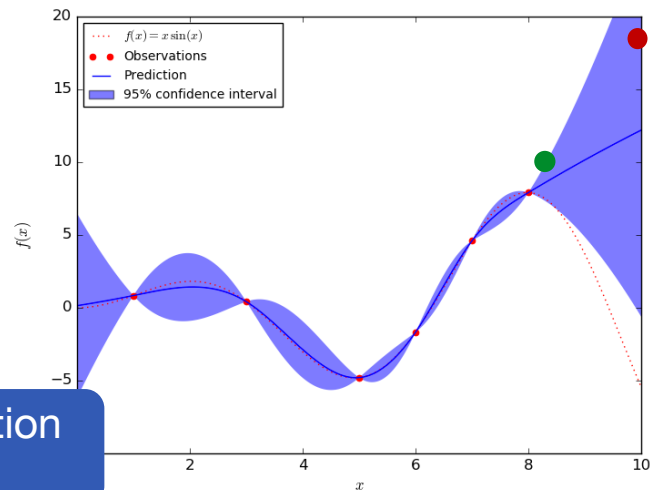    $S \leftarrow S \cup \{v_t\}$
    Receive feedback $y_t = f(v_t) + \epsilon_t$
**end for**

Learn posterior

Trades off exploration exploitation
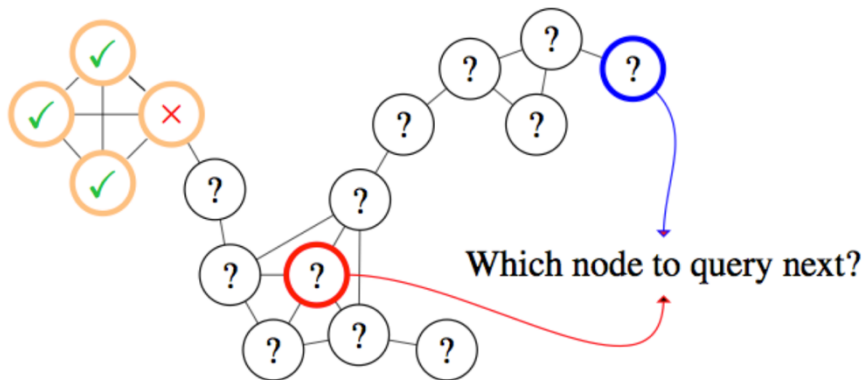
Ask user feedback



- **Exploration**: select items with high-variance
- **Exploitation**: select items with high-value

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Active learning on graphs – which prior?

[Ma et al., 2015]

**Idea:** Use the graph structure to infer the node classes

Use graph Laplacian as prior
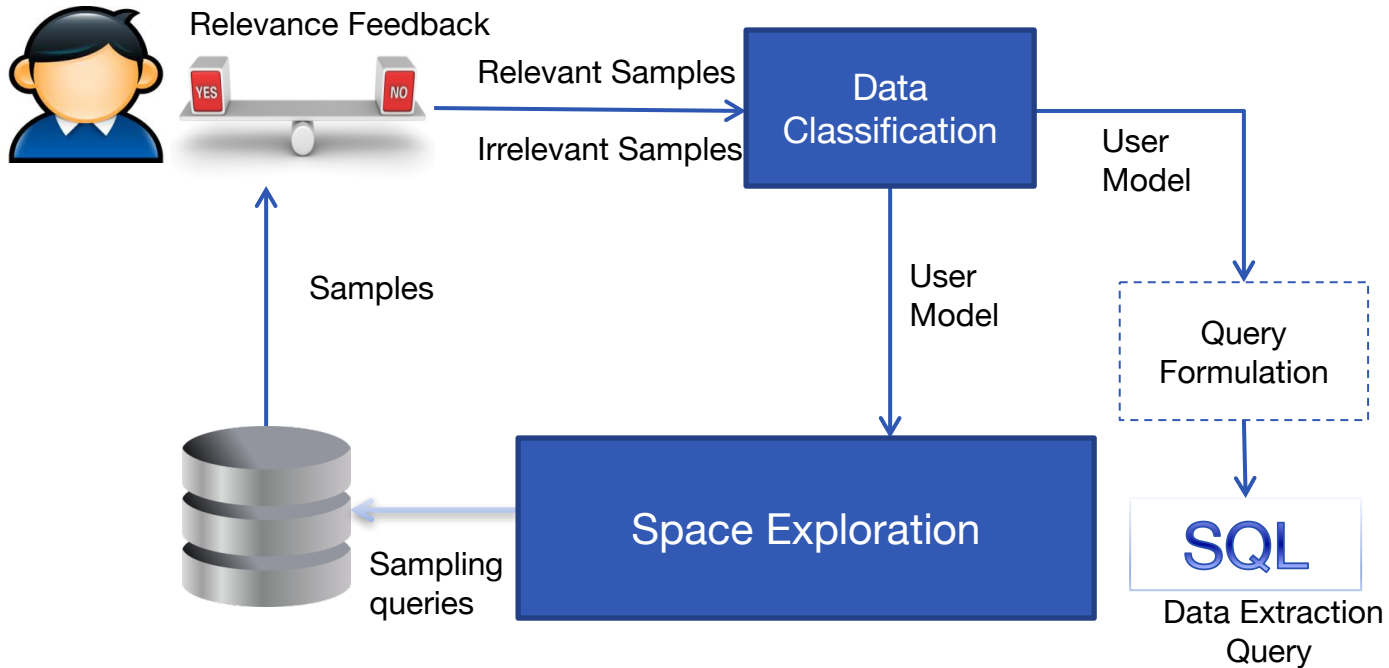$L = D - A$, A is the adjacency matrix



Which node to query next?

$$p(\mathbf{f}) \sim \mathcal{N}(0, L^{-1})$$

Laplacian: higher probability of having the same class if two nodes are connected

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Explore-by-Example: AIDE

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# The AIDE algorithm
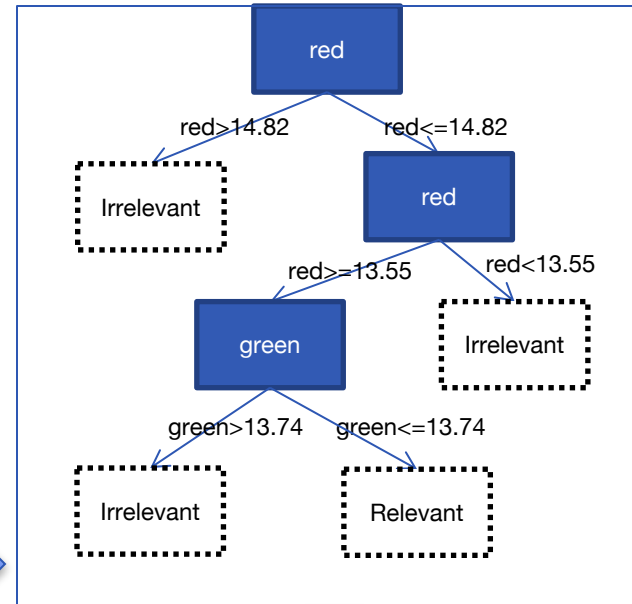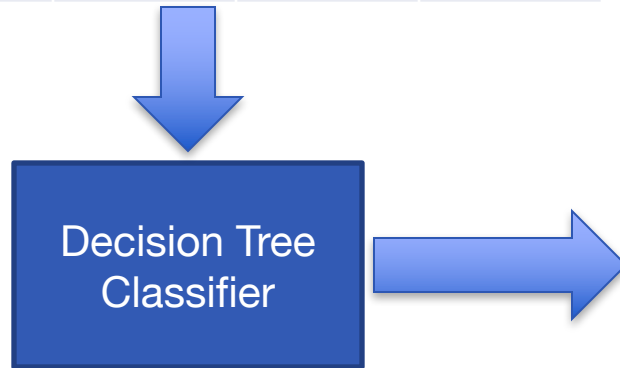
1. Divide the space into d-dimensional cubes
2. Find the sample points in the cubes (medoids)
3. Train the classifier
4. Refine the training sampling from neighbors of misclassified points
5. Boundary refinement

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Classification & Query Formulation

[Dimitriadou et al., 2015]

| Sample | Red | Green | Relevant |
|--------|-------|-------|----------|
| Object A | 13.67 | 12.34 | Yes |
| Object B | 15.32 | 14.50 | No |
| .. | .. | .. | ... |
| Object X | 14.21 | 13.57 | Yes |

**Decision Tree Classifier**

red

red>14.82 → Irrelevant

red<=14.82 → red

red>=13.55 → green

red<13.55 → Irrelevant

green>13.74 → Irrelevant

green<=13.74 → Relevant

SELECT * FROM galaxy WHERE red<= 14.82 AND red>= 13.5 AND green<=13.74

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI db Trento

# Misclassified Sample Exploitation

Sampling Areas

Red wavelength

Green Wavelength

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Clustering-based Sampling

**Idea**: Use a k-medoid approach to find sampling areas

Clusters-Sampling Areas

Red wavelength

Green Wavelength

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Where we are

Relational databases

Textual data

Graphs and networks

Machine learning

Challenges and Remarks

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI    db Trento

# Example-based methods

- Query suggestion using examples
- Reverse engineering queries



- Entity extraction by example text
- Web table completion using examples
- Search by example



- Community-based Node-retrieval
- Entity Search
- Path and SPARQL queries
- Graph structures as Examples

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Example-based methods: takeaways

## Relational

- Complex search space
- Exact and approximate
- Interactivity can improve the quality
- Limited to query inference

## Textual

- Allows serendipitous search
- Easier document finding
- Speed up entity matching

## Graph

- Exploit locality
- Entity attributes are expressive
- Reverse engineering: good approximations
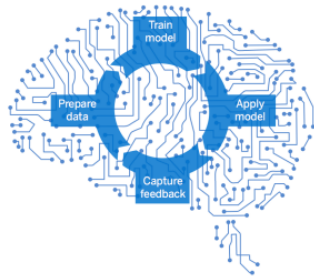- Large result-sets require ranking

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# The use of examples

## Examples can ease data exploration

- … reduce need for complex queries / simplify user input
- … require no schema knowledge
- … allow uncertainity in search conditions
- … require little data analytics expertise

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Where should we invest time



**Machine learning**

**Approximate Methods**

**User models**

**Scalability**

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI ᴅᑲ Trento

# ADOPT HETEROGENEITY

**Need for solutions that**
operate across different models

operate on heterogeneous datastores

"The Context of Mobile Interaction"
– Nadav Savio

Nadav Savio | Giant Ant Design | www.giantant.com

# PERSONALIZATION

**better understand user needs**

**Meta-data and User Profiles**

exploit *query log, prior searches, user context*

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI · dbTrento

Design: M. Lissandrini

# **DEMOCRATIZATION**
## **easy access to data**

tools that work on
commodity
hardware, mobile
devices

data-exploration for
everyday use-cases

# INTERACTIVITY

gradually understand
user need

# ADAPTIVITY

build indexes and data
structures on-the-go

D. Mottin, M. Lissandrini

# NATARUAL LANGUAGE INTERFACE

flexible, vague, imprecise input

exploration through conversation

D. Mottin, M. Lissandrini

HPI db Trento

# Example is always more efficacious than precept

**Samuel Johnson**, *Rasselas (1759), Chapter 29.*

**"New Trends on Exploratory Methods for Data Analytics."**

Davide Mottin, Matteo Lissandrini, Yannis Velegrakis, Themis Palpanas.

*Proceedings of the Conference in Very Large Databases (PVLDB), 10(12), 2017*

**Slides:** http://j.mp/DataExplore

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

HPI dbTrento

# Acknowledgments

**We would like to thank the authors of the papers who kindly provided us the slides**

Angela Bonifati, Radu Ciucianu, Marcelo Arenas, Gonzalo Diaz, Egor Kostylev, Yaacov Weiss, Sarah Cohen, Fotis Psallidas, Li Hao, Chan Chee Yong, Ilaria Bordino, Mohamed Yakout, Kris Ganjam, Kaushik Chakrabati, Thibault Sellam, Rohit Singh, Maeda Hanafi, Marcin Sydow, Mingzhu Zhu, Yoshiharu Ishikawa, Daniel Deutch, Nandish Jayaram, Bryan Perozzi, Kiriaki Dimitriadou, Yifei Ma, Natali Ruchansky, Quoc Trung Tran, Hastagiri Prakash Vanchinathan

# References

M. Arenas, G. I. Diaz, and E. V. Kostylev**. Reverse engineering sparql queries**. WWW, 2016.

Agichtein, E. and Gravano, L. **Snowball: Extracting relations from large plain-text collections.** ICDL, 2000.

A.Bonifati, R.Ciucanu,and A.Lemay. **Learning path queries on graph databases**. EDBT, 2015.

A. Bonifati, R. Ciucanu, and S. Staworko. **Learning join queries from user examples**. TODS, 2016.

A. Bonifati, U. Comignani, E. Coquery, and R. Thion. **Interactive mapping specification with exemplar tuples.** SIGMOD, 2017.

I. Bordino, G. De Francisci Morales, I. Weber, and F. Bonchi. **From machu picchu to rafting the urubamba river: anticipating information needs via the entity-query graph**. WSDM, 2013.

D. Deutch and A. Gilad. **Qplain: Query by explanation**. ICDE, 2016.

D. Mottin, M. Lissandrini

# References

G. Diaz, M. Arenas, and M. Benedikt. **Sparqlbye: Querying rdf data by example**. PVLDB, 2016.

K. Dimitriadou, O. Papaemmanouil, and Y. Diao. **Explore-by-example: An automatic query steering framework for interactive data exploration.** In SIGMOD, 2014.

B. Eravci and H. Ferhatosmanoglu. **Diversity based relevance feedback for time series search.** PVLDB, 2013.

A. Gionis, M. Mathioudakis, and A. Ukkonen. **Bump hunting in the dark: Local discrepancy maximization on graphs**. ICDE, 2015.

M. F. Hanafi, A. Abouzied, L. Chiticariu, and Y. Li. **Synthesizing extraction rules from user examples with seer.** SIGMOD, 2017.

Y. Ishikawa, R. Subramanya, and C. Faloutsos. **Mindreader: Querying databases through multiple examples.** VLDB, 1998.

N. Jayaram, A. Khan, C. Li, X. Yan, and R. Elmasri. **Querying knowledge graphs by example entity tuples**. TKDE, 2015.

D. Mottin, M. Lissandrini

HPI db Trento

# References

H. Li, C.-Y. Chan, and D. Maier. **Query from examples: An iterative, data-driven approach to query construction.** PVLDB, 2015.

Y. Ma, T.-K. Huang, and J. G. Schneider. **Active search and bandits on graphs using sigma-optimality.** UAI, 2015.

S. Metzger, R. Schenkel, and M. Sydow. **Qbees: query by entity examples.** CIKM, 2013.

D. Mottin, M. Lissandrini, Y. Velegrakis, and T. Palpanas. **Searching with xq: the exemplar query search engine**. SIGMOD, 2014.

D. Mottin, M. Lissandrini, Y. Velegrakis, and T. Palpanas. **Exemplar queries: a new way of searching.** VLDB J., 2016.

B. Perozzi, L. Akoglu, P. Iglesias Sa´nchez, and E. Müller. **Focused clustering and outlier detection in large attributed graphs.** KDD, 2014.

# References

F. Psallidas, B. Ding, K. Chakrabarti, and S. Chaudhuri**. S4: Top-k spreadsheet-style search for query discovery**. SIGMOD, 2015.

R. Rolim, G. Soares, L. D'Antoni, O. Polozov, S. Gulwani, R. Gheyi, R. Suzuki, and B. Hartmann. **Learning syntactic program transformations from examples**. ICSE, 2017.

N. Ruchansky, F. Bonchi, D. García-Soriano, F. Gullo, and N. Kourtellis. **The minimum wiener connector problem**. SIGMOD, 2015.

T. Sellam and M. Kersten. **Cluster-driven navigation of the query space.** TKDE, 2016.

Y. Shen, K. Chakrabarti, S. Chaudhuri, B. Ding, and L. Novik. **Discovering queries based on example tuples.** SIGMOD, 2014.

R. Singh. **Blinkfill: Semi-supervised programming by example for syntactic string transformations.** PVLDB, 2016.

G. Sobczak, M. Chochół, R. Schenkel, and M. Sydow. iqbees: **Towards interactive semantic entity search based on maximal aspects**. Foundations of Intelligent Systems, 2015.

D. Mottin, M. Lissandrini

HPI db Trento

# References

Y. Su, S. Yang, H. Sun, M. Srivatsa, S. Kase, M. Vanni, and X. Yan. **Exploiting relevance feedback in knowledge graph search**. KDD, 2015.

Q. T. Tran, C.-Y. Chan, and S. Parthasarathy. **Query reverse engineering**. VLDB J., 2014.

H. P. Vanchinathan, A. Marfurt, C.-A. Robelin, D. Koss- mann, and A. Krause. **Discovering valuable items from massive data**. In KDD, 2015.

C.Wang, A.Cheung, and R.Bodik**. Interactive query synthesis from input-output examples.** In SIGMOD, 2017.

C. Wang, A. Cheung, and R. Bodik. **Synthesizing highly expressive sql queries from input-output examples.** In PLDI, 2017.

Y. Y. Weiss and S. Cohen. **Reverse engineering spj-queries from examples**. SIGMOD, 2017.

M. Yakout, K. Ganjam, K. Chakrabarti, and S. Chaudhuri. **Infogather: Entity augmentation and attribute discovery by holistic matching with web tables.** SIGMOD, 2012.

M. Zhu and Y.-F. B. Wu. **Search by multiple examples**. WSDM, 2014.

M. M. Zloof. **Query by example**. AFIPS NCC, 1975.