# Example-driven Search:
## a New Frontier for Exploratory Search

Matteo Lissandrini (Aalborg University),       Yannis Velegrakis (Utrecht University)

Davide Mottin (Aarhus University),       Themis Palpanas (University of Paris)

AARHUS UNIVERSITY

AALBORG UNIVERSITET

Utrecht University

UNIVERSITÉ PARIS DESCARTES

**Link for questions**

https://j.mp/ExploreSIGIR

**Tutorial Slides and Other Material**

https://data-exploration.ml/

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

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

# Traditional Data Management Systems



rdf:Type=GlutamateReceptor
ro:has_function=ex:GlutamateReceptorActivity001

Data

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

# Modern Data Management Systems

## Not clear what we are looking for

I would like to find
acquisitions
like the one of
YouTube by Google

Big Data

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

# Exploration

*We know where we start
we don't know what we'll find*

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

# Exploration



Traditional

On data

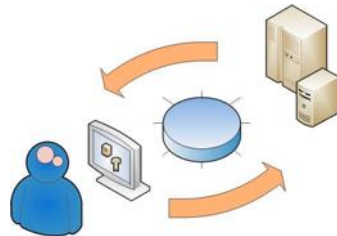M. Lissandrini, D. Mottin, T. Palpanas, Y. Velegrakis

# Data exploration



Cleaning and profiling



Visualization



Analysis



Interactions



Architectures

M. Lissandrini, D. Mottin, T. Palpanas, Y. Velegrakis
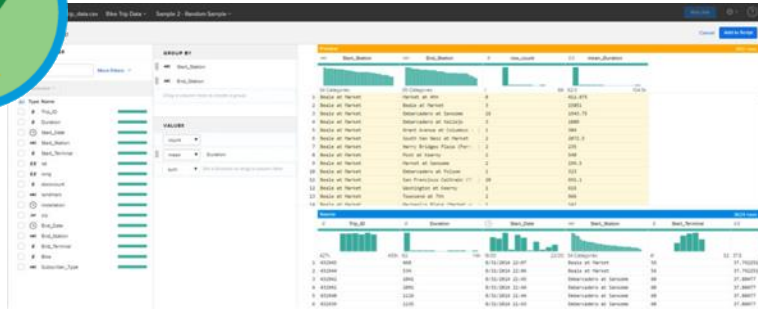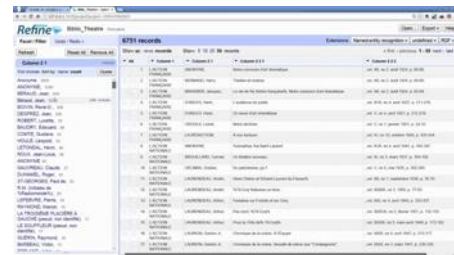
# Data exploration software


Trifacta: data preparation


OpenRefine: data preparation and cleanup


Tableau: analysis and statistics

M. Lissandrini, D. Mottin, 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

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Modern Data Management Systems

How do we describe what we are looking for?

M. Lissandrini, D. Mottin, 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

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

# Examples as Exploratory Methods

**Example is always more efficacious than precept**

Samuel Johnson, Rasselas (1759), Chapter 29.

Is there a galaxy like this?

Answers

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Tutorial's goals

**Techniques, Algorithms, Applications for using Examples to support Exploratory**

- 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

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

# Our book on Example-based methods



**Matteo Lissandrini**
*Aalborg University*

Knowledge Graphs , Novel Query Paradigms,
Graph Mining
http://people.cs.aau.dk/~matteo

**Yannis Velegrakis**
*Utrecht University*

Big Data Management & Analytics, Information
Integration, Data Curation
https://velgias.github.io

**Davide Mottin**
*Aarhus University*

Graph Mining, Novel Query Paradigms,
Interactive Methods
https://mott.in

**Themis Palpanas**
*Paris Descartes University*

Data Series Indexing & Mining, Data Analytics &
Management
http://www.mi.parisdescartes.fr/~themisp

MORGAN & CLAYPOOL PUBLISHERS

# Data Exploration Using Example-Based Methods

Matteo Lissandrini
Davide Mottin
Themis Palpanas
Yannis Velegrakis

SYNTHESIS LECTURES ON DATA MANAGEMENT

H. V. Jagadish, Series Editor

**SIGIR 2019** tutorial
https://data-exploration.ml

M. Lissandrini, D. Mottin, 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    | *          | *     | *    |

Incomplete values

Unspecified values

Value conditions

| Intuitive interface for simple queries | Restricted to SQL syntax but not explicitly |
|---|---|
| SQL not required | Not example-based |

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

# Similarities are the key …

If we knew how similar each item is with respect to any other for **each** user, we would know the answer to

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

# Similarities are the key ...

**We define:**

A universe $\mathcal{U}$ of items

A similarity among items ~

A set of input examples $\mathcal{E}$

A set of output user desired answers $\mathcal{A}$



Universe #

Desired Answers $

Examples $\mathcal{E}$

Similarity relation ~

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

# The example-based problem

**Given**

   a set of examples $\mathcal{E}$ from a universe $\mathcal{U}$

**Find**

<div align="center">a similarity $\sim$ such that</div>

1.   $\mathcal{E}$ is part of the answers $\mathcal{A}$ partially or totally

2.   The answers in $\mathcal{A}$ are the most similar to the examples in $\mathcal{E}$ according to $\sim$

<div align="center">How do we find $\sim$ for each user?<br>
Do we need to know exactly $\sim$?</div>

**SIGIR 2019** tutorial
https://data-exploration.ml                    M. Lissandrini, D. Mottin, T. Palpanas, Y. Velegrakis

# Example-based methods

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

# Example-based methods

## Relational

Reverse engineering queries

Example-driven schema mapping

Interactive data repairing

## Textual

Entity extraction by example text

Web table completion using examples

Search by example

## Graph

Community-based Node-retrieval

Entity Search

Path and SPARQL queries Graph structures as Examples

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

# Tutorial structure

Relational databases

Textual data

Graph and networks

Machine learning

Challenges and Remarks

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Where we are

Relational databases

Textual data

Graphs and networks

Machine learning

Challenges and Remarks

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Searching for …

**SEARCHING** FOR

**Tuples**

Matching rules

BY **FOCUSING ON**

Exact Queries

Approximate Queries

Schemas

Data cleaning

**APPLYING**

**One-step**
[Tran et al. 2009,
Zhang et al., 2013,
Weiss and Cohen, 2017]

**Interactive**
[Li et al., 2015a]

**Minimal**
[Shen et al., 2014
Tran et al., 2014]

**Top-k**
[Psallidas et al., 2015
Panev and Michel, '16]

**Schema to examples**
[Alexe et al., 2011a]

**Example-driven**
[Alexe et al. 2011b,
Cate et al. 2013,
Gottlob Senellart, 2010
Bonifati et al., 2016]

**Entity Mentions**
[Singh et al. 2017]

**Data repairing**
[He et al., 2016]

**PRODUCES**

Reverse Engineered SQL queries

Data Integration

Duplicate matches

M. Lissandrini, D. Mottin, T. Palpanas, Y. Velegrakis
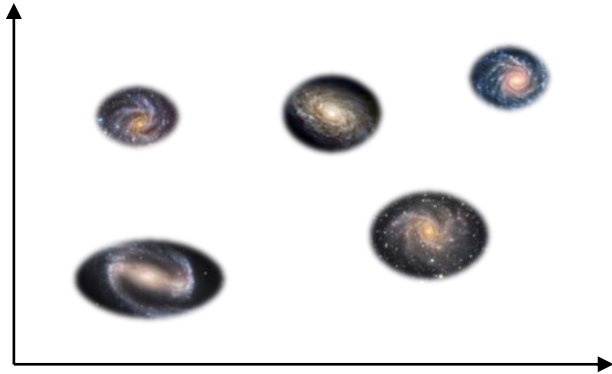
# Reverse engineering queries (REQ)

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

Example tuples

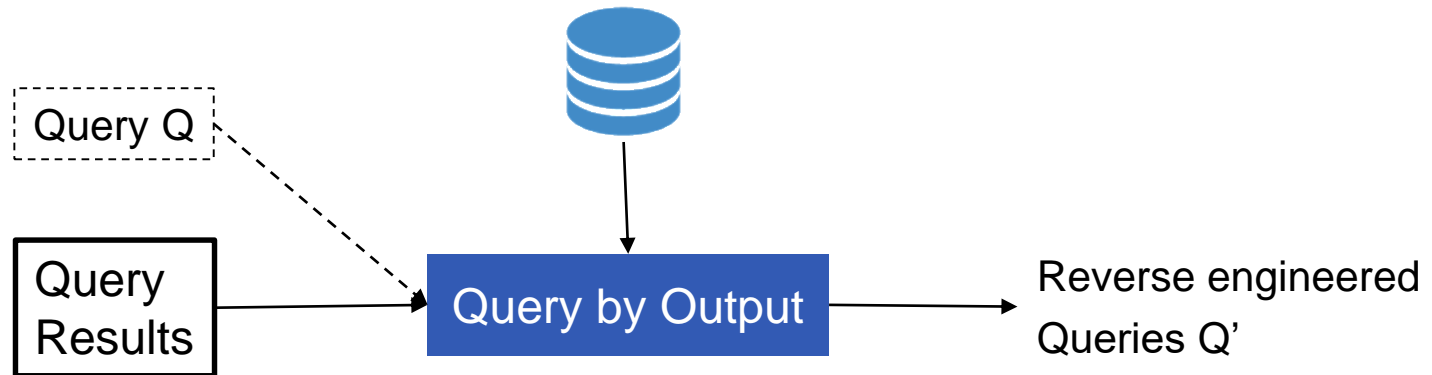How do you find such queries?

```
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
```

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

# Query by Output – TALOS (classification-based)

**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'

# How many reverse engineered queries?

Master

| | name | bat | throw | stint | weight | 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 |

| | | |
|---|---|---|
| $r_1$ | B | PIT |
| $r_2$ | E | CHA |

$Q(D)$

## What queries generated Q(D)?

Q1 = SELECT name, team FROM Master WHERE bat = 'R' AND throw = 'R'

Q2 = SELECT name, team FROM Master WHERE bat = 'R' AND weight > 35

Q3 = SELECT name, team FROM Master WHERE bat = 'R' AND stint <> 2
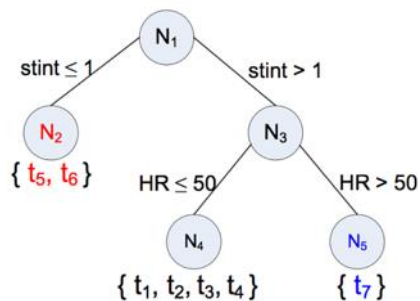
...

Instance Equivalent Queries

# TALOS

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

|       | name | bat | throw | stint | HR | team |   |
|-------|------|-----|-------|-------|-----|------|---|
| $t_1$ | A    | L   | R     | 2     | 40  | PIT  | X |
| $t_2$ | A    | L   | R     | 2     | 50  | MT1  | X |
| $t_3$ | C    | R   | L     | 2     | 35  | CHA  | X |
| $t_4$ | D    | L   | R     | 3     | 30  | PIT  | X |
| $t_5$ | 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

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

# How complex is exact REQ?

[Weiss et al., 2017]

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?

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Complexity  - No parameters

[Weiss et al., 2017]

| 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
Bounded selections: HARD

Combination of operators: HARD!!!

Reduction from SAT

**SIGIR 2019** tutorial
https://data-exploration.ml

M. Lissandrini, D. Mottin, 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   $D \neq 3$

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

# Bounded select

INPUT: Database D, Examples E, Query size k

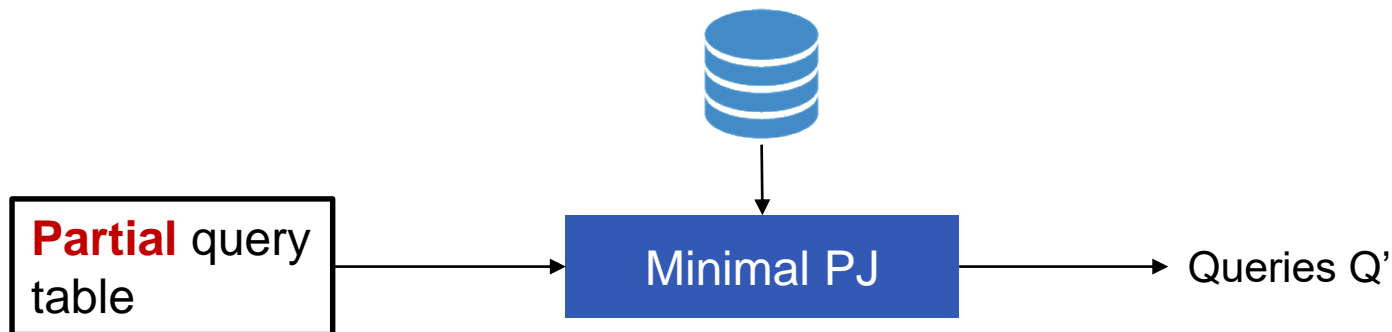OUTPUT: Does there exist a query satisfying D and E, of size at most k?

$$U = \{1,2,3,4,5\} \qquad S = \{ \{1,2,3\}, \{2,4\}, \{3,4\}, \{4,5\} \}$$

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

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Minimal Project Join REQ

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

| **Partial** query table | → | Minimal PJ | → | Queries Q' |
| --- | --- | --- | --- | --- |

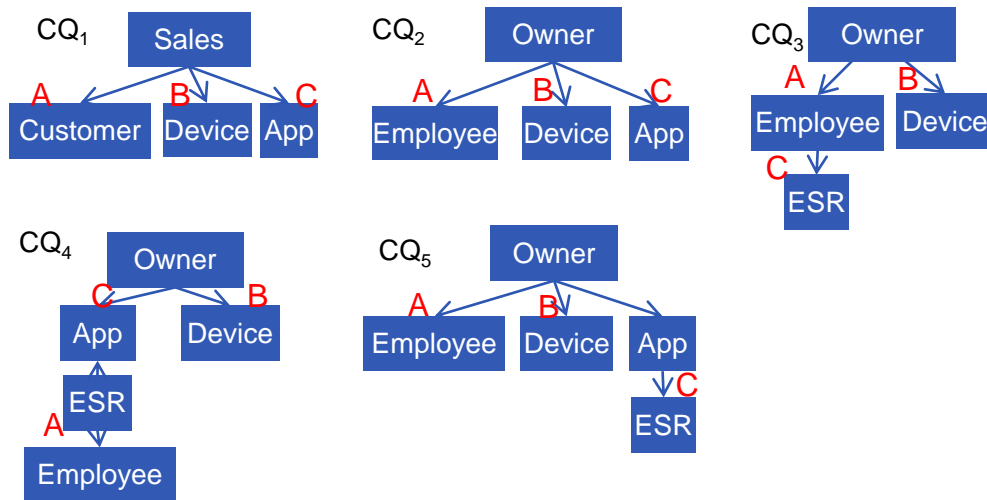|   | A | B | C |
| --- | --- | --- | --- |
| 1 | Mike | ThinkPad | Office |
| 2 | Mary | iPad |  |
| 3 | Bob |  | Dropbox |

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

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

# Candidate Query Generation

[Shen et al., 2014]

- Use candidate network generation algorithm (Hristidis 2002)

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



$CQ_1$ — Sales → A Customer, B Device, C App

$CQ_2$ — Owner → A Employee, B Device, C App

$CQ_3$ — Owner → A Employee, B Device, C ESR

$CQ_4$ — Owner → C App → A Employee (ESR), B Device

$CQ_5$ — Owner → A Employee, B Device, App → C ESR

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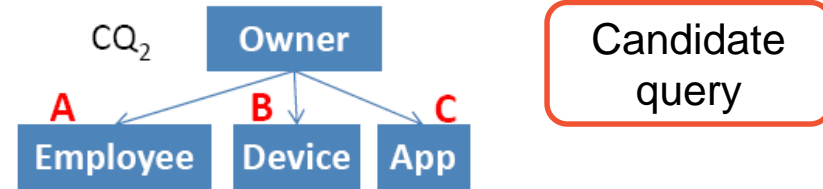   - Every leaf node contains a column that is mapped by an input column

# Validity verification

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

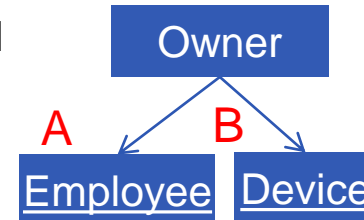Better: exploit substructures in candidate queries for pruning

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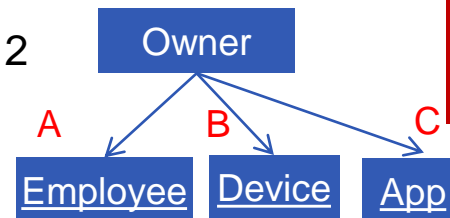
NP-hard



Candidate query

Substructures

Sub 1 fails => $CQ_2$ invalid

Sub 1 fails => Sub 2 fails

# 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 |

# Ranking score
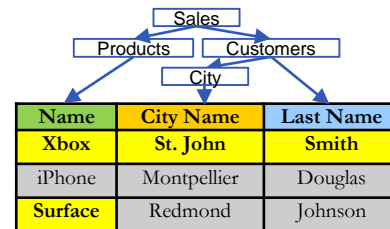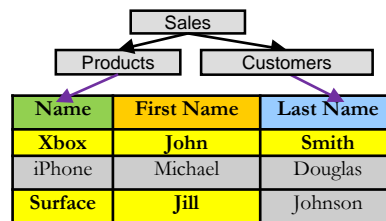
Linear combination of row score and column score

(Overlapping with the example table)

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

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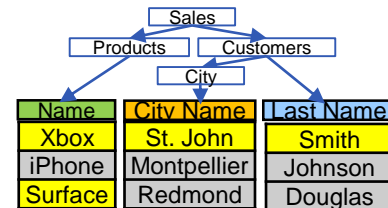- $\alpha = 0$ penalizes missing columns

**Row score**

| | | | Row Score | |
|---|---|---|---|---|
| John | Smith | Xbox | 3 | 3 |
| Jill | Hans | Surface | 2 | 1 |
| | | | 5 | 4 |

**Column score**

| | John | Smith | Xbox | |
|---|---|---|---|---|
| | Jill | Hans | Surface | |
| Column Score | 2 | 1 | 2 | 5 |
| | 2 | 1 | 1 | 4 |

Sales → Products, Customers

| Name | First Name | Last Name |
|---|---|---|
| Xbox | John | Smith |
| iPhone | Michael | Douglas |
| Surface | Jill | Johnson |

Sales → Products, Customers → City

| Name | City Name | Last Name |
|---|---|---|
| Xbox | St. John | Smith |
| iPhone | Montpellier | Douglas |
| Surface | Redmond | Johnson |

Sales → Products, Customers

| Name | First Name | Last Name |
|---|---|---|
| Xbox | John | Smith |
| iPhone | Jill | Johnson |
| Surface | Michael | Douglas |

Sales → Products, Customers → City

| Name | City Name | Last Name |
|---|---|---|
| Xbox | St. John | Smith |
| iPhone | Montpellier | Johnson |
| Surface | Redmond | Douglas |

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Interactive REQ – Query from Examples (cost model)

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

Modified database and results

Query Results

REQ

Database Refinement

Reverse engineered Queries Q'

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

# Database Refinement

[Li et al., 2015]

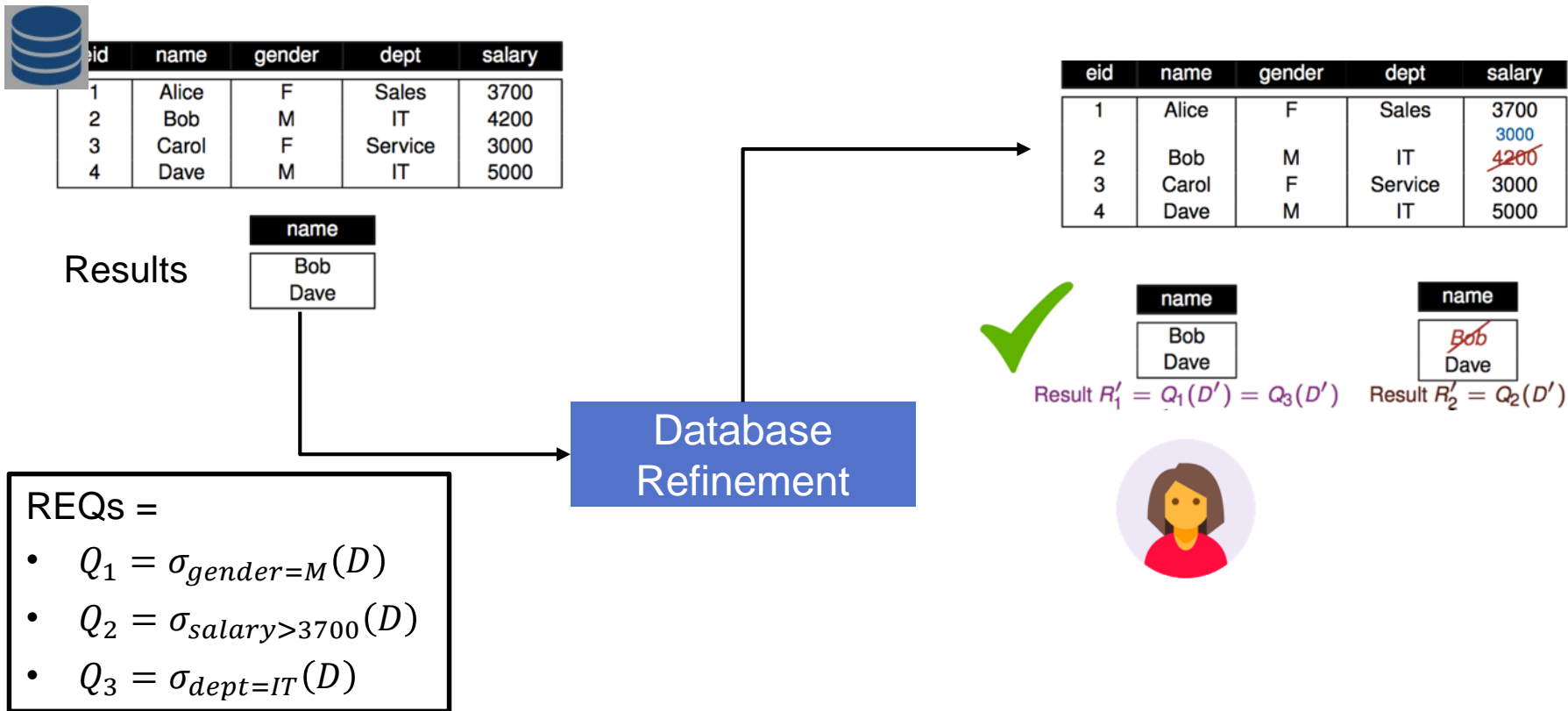| eid | name | gender | dept | salary |
|-----|------|--------|------|--------|
| 1 | Alice | F | Sales | 3700 |
| 2 | Bob | M | IT | 4200 |
| 3 | Carol | F | Service | 3000 |
| 4 | Dave | M | IT | 5000 |

Results

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

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

**Database Refinement**

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

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

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

Result $R_2' = Q_2(D')$

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

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Cost model

Number of modified tables

Choice to go with the max tuples modified

Number of new result sets

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

DB cost   Results cost   Effort to examine D'   Effort to examine new results
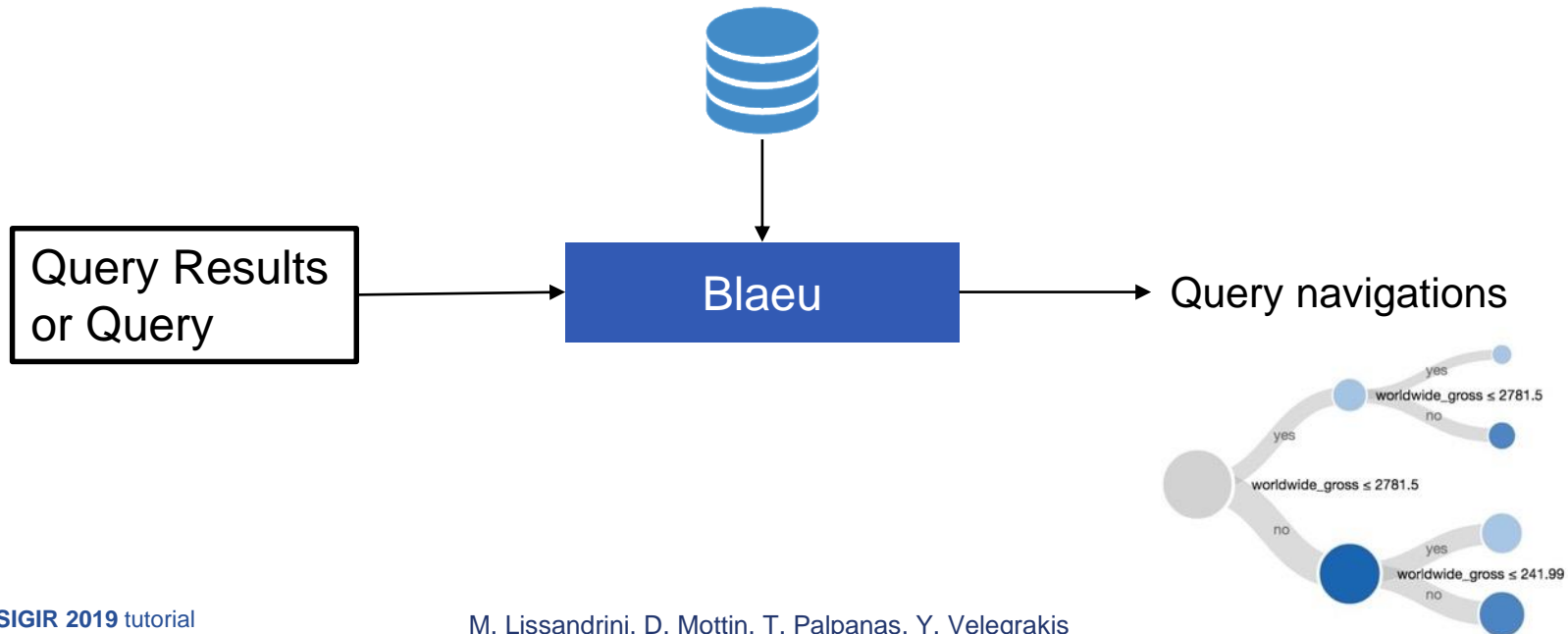
Current cost   Residual cost

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
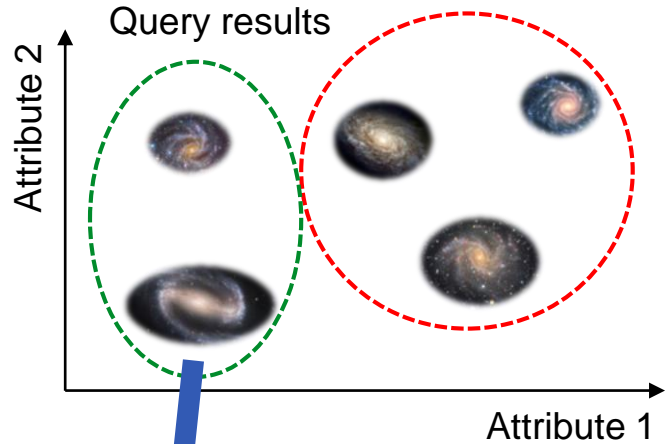3. The query are balanced (less interactions)

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

# Examples for query suggestion: Blaeu (Clustering)  [Sellam et al., 2016]

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



Query Results or Query → Blaeu → Query navigations

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Examples for query suggestion: Blaeu    [Sellam et al., 2016]

Query results



Attribute 2

Attribute 1

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

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

User utility

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

**SIGIR 2019** tutorial
https://data-exploration.ml

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

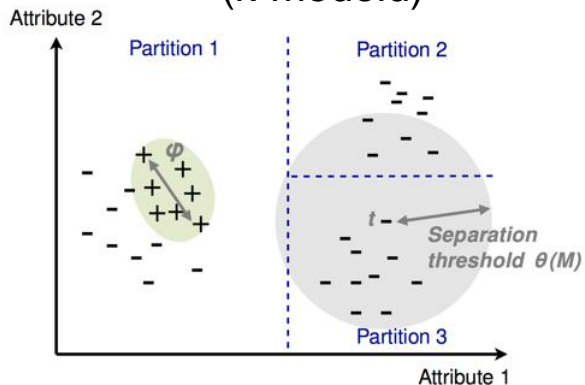# Examples for query suggestion: Blaeu [Sellam et al., 2016]

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

Unknown User utility

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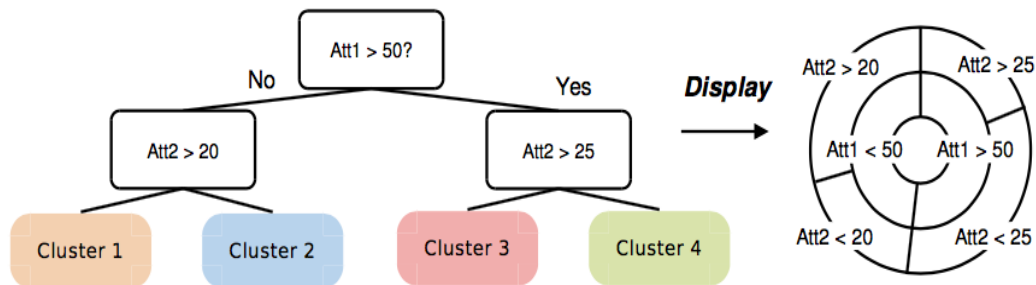
**Solution**: interesting tuples are close to each other within a maximum separation threshold $\theta(\mathcal{C})$
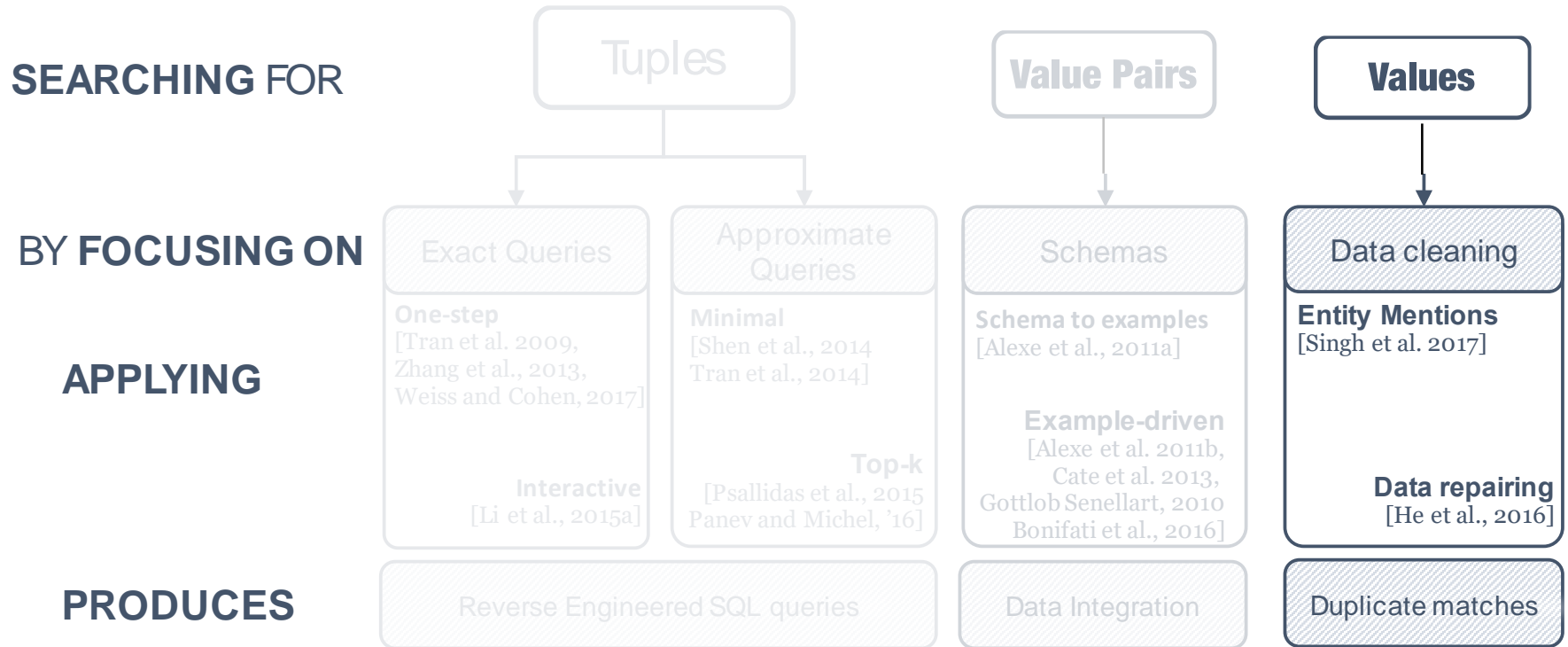
Detect clusters
(k-medoid)

Organize clusters

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Searching for …

| | Tuples | | Value Pairs | Values |
|---|---|---|---|---|
| **SEARCHING** FOR | | | | |
| BY **FOCUSING ON** | Exact Queries | Approximate Queries | Schemas | Data cleaning |
| **APPLYING** | One-step [Tran et al. 2009, Zhang et al., 2013, Weiss and Cohen, 2017] | Minimal [Shen et al., 2014 Tran et al., 2014] | Schema to examples [Alexe et al., 2011a] | **Entity Mentions** [Singh et al. 2017] |
| | Interactive [Li et al., 2015a] | Top-k [Psallidas et al., 2015 Panev and Michel, '16] | Example-driven [Alexe et al. 2011b, Cate et al. 2013, Gottlob Senellart, 2010 Bonifati et al., 2016] | Data repairing [He et al., 2016] |
| **PRODUCES** | Reverse Engineered SQL queries | | Data Integration | Duplicate matches |

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

# Data Cleaning

- Often data have redundancy, wrong values, and missing values

- Different values can represent the same object (e.g., N.Y. and New York)

- Values can be simply wrong

Data cleaning refers to ways of making the data consistent and correct

| tid | Date | Molecule | Laboratory | Quantity |
|-----|------|----------|------------|----------|
| t1 | 11 Nov | $C_{16}H_{16}Cl$ | Austin | 200 |
| t2 | 12 Nov | statin | Austin | 100 |
| t3 | 12 Nov | $C_{24}H_{75}S_6$ | N.Y. | 100 |
| t4 | 12 Nov | statin | Boston | 200 |
| t5 | 13 Nov | statin | Austin | 200 |
| t6 | 15 Nov | $C_{17}H_{20}N$ | Dubai | 1000 |

| tid | Date | Molecule | Laboratory | Quantity |
|-----|------|----------|------------|----------|
| t1 | 11 Nov | $C_{16}H_{16}Cl$ | Austin | 200 |
| t2 | 12 Nov | $C_{22}H_{28}F$ | Austin | 100 |
| t3 | 12 Nov | $C_{24}H_{75}S_6$ | New York | 100 |
| t4 | 12 Nov | statin | Boston | 200 |
| t5 | 13 Nov | $C_{22}H_{28}F$ | Austin | 200 |
| t6 | 15 Nov | $C_{17}H_{20}N$ | Dubai | 100 |

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Data repairing: rules

[He, J. et al. 2016]

A **rule** is a logical formula which determines how to change the value in a cell or a group of cells.

IF $[X_1 = C_1 \dots X_n = C_n]$ UPDATE $X_i$ to some value

- The update $t_3[\text{Laboratory}] \leftarrow$ "New York" can be obtained by the rule

- IF [Laboratory = "N.Y."] UPDATE Laboratory to "New York"

- ```
  UPDATE Table
  SET Laboratory='New York'
  WHERE tid=t3
  ```

   BUT it needs to be done for each cell!!

| tid | Date | Molecule | Laboratory | Quantity |
|-----|------|----------|------------|----------|
| t1 | 11 Nov | $C_{16}H_{16}Cl$ | Austin | 200 |
| t2 | 12 Nov | statin | Austin | 100 |
| t3 | 12 Nov | $C_{24}H_{75}S_6$ | N.Y. | 100 |
| t4 | 12 Nov | statin | Boston | 200 |
| t5 | 13 Nov | statin | Austin | 200 |
| t6 | 15 Nov | $C_{17}H_{20}N$ | Dubai | 1000 |

| tid | Date | Molecule | Laboratory | Quantity |
|-----|------|----------|------------|----------|
| t1 | 11 Nov | $C_{16}H_{16}Cl$ | Austin | 200 |
| t2 | 12 Nov | $C_{22}H_{28}F$ | Austin | 100 |
| t3 | 12 Nov | $C_{24}H_{75}S_6$ | New York | 100 |
| t4 | 12 Nov | statin | Boston | 200 |
| t5 | 13 Nov | $C_{22}H_{28}F$ | Austin | 200 |
| t6 | 15 Nov | $C_{17}H_{20}N$ | Dubai | 100 |

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

**UPDATES:**

$\Delta_1$: t3[Laboratory] ← "New York"
$\Delta_2$: t6[Quantity] ← 100
$\Delta_3$: t2[Molecule] ← "$C_{22}H_{28}F$"

**Some rules for $\Delta_1$:**
1. Change all Laboratory values to "New York" (t1 – t6)
2. Reformatting all "N.Y" to "New York"(t3)

**Some rules for $\Delta_2$:**
1. Update the quantity to 100 if the molecule is $C_{17}H_{20}N$ and the date is 15 Nov (t6)

**Some rules for $\Delta_3$:**
1. Update to "$C_{22}H_{28}F$" if molecule is statin (t2,t4,t5)
2. Update to "$C_{22}H_{28}F$" if molecule is statin and Laboratory Austin (t2,t5)
3. Update to "$C_{22}H_{28}F$" if molecule is statin and lab is Austin and date is 12 Nov and quantity is 100 (t2)

| tid | Date | Molecule | Laboratory | Quantity |
|---|---|---|---|---|
| t1 | 11 Nov | $C_{16}H_{16}Cl$ | Austin | 200 |
| t2 | 12 Nov | $C_{22}H_{28}F$ | Austin | 100 |
| t3 | 12 Nov | $C_{24}H_{75}S_6$ | New York | 100 |
| t4 | 12 Nov | statin | Boston | 200 |
| t5 | 13 Nov | $C_{22}H_{28}F$ | Austin | 200 |
| t6 | 15 Nov | $C_{17}H_{20}N$ | Dubai | 100 |

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Interactive data cleaning: problem

User validates rules, but has no capacity to validate all rules for each update.

- **Budget Repair Problem:** Given a set $Q$ of rules, a table T and a budget B, **find B rules from $Q$ to maximize the number of repairs over T**

- Budget repair problem is an *online problem*

Corresponding *offline problem* is: given as input $Q$ rules where validity of each rule is known, select B rules from $Q$ to maximize the number of repairs over T. (NP-Hard)

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

# Rule lattice

**More specific**

| tid | Date | Molecule | Laboratory | Quantity |
|---|---|---|---|---|
| t2 | 12 Nov | statin→$C_{22}H_{28}F$ | Austin | 100 |

DMLQ(1)

→ Number of tuples affected

DML (1)    DMQ (1)    DLQ (1)    MLQ (1)

DM (1)    DL (1)    DQ (2)    ML (2)    MQ (1)    LQ (1)

D (3)    M (3)    L (3)    Q (2)

**12 Nov**    **statin**    **Austin**    **100**

∅ (6)

**More general**

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

# Lattice pruning

1. If Q is valid, Q' is also valid if Q' $\leqslant$ Q
2. If Q is invalid, Q'' is also invalid if Q $\leqslant$ Q''
3. If Q is valid, all Q' such that Q' $\leqslant$ Q are valid.
4. If Q is invalid, all Q'' such that  Q $\leqslant$ Q'' are invalid.

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Lattice pruning



**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Dive search

- Binary Search over the lattice ,ordering with #affected tuples
- If T → BinarySearch(Q$_\wedge$)
- If F → BinarySearch(Q $_\vee$)

BinarySearch()

Q1

Q2

BinarySearch()

BinarySearch()

## Algorithm complexity
$$\mathcal{O}(B|\mathcal{Q}|\log|\mathcal{Q}|)$$

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

# Searching for …

| | Tuples | | Value Pairs | Values |
|---|---|---|---|---|
| **SEARCHING** FOR | | | | |
| BY **FOCUSING ON** | Exact Queries | Approximate Queries | Schemas | Data cleaning |
| **APPLYING** | One-step [Tran et al. 2009, Zhang et al., 2013, Weiss and Cohen, 2017] | Minimal [Shen et al., 2014 Tran et al., 2014] | **Schema to examples** [Alexe et al., 2011a] | Entity Mentions [Singh et al. 2017] |
| | Interactive [Li et al., 2015a] | Top-k [Psallidas et al., 2015 Panev and Michel, '16] | **Example-driven** [Alexe et al. 2011b, Cate et al. 2013, Gottlob Senellart, 2010 Bonifati et al., 2016] | Data repairing [He et al., 2016] |
| **PRODUCES** | Reverse Engineered SQL queries | | Data Integration | Duplicate matches |

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

# Schema mapping

- Schema mapping finds a way to represent items on one database to items on another database
- Finds a mapping Σ between two schemas such that a query on one database can be converted to a query on the other database
- Schema mappings in Σ are rules in first-order logic that specifies the relationships between schema S and T

$$\forall x \forall y \, S(x, y) \wedge U(x, z) \rightarrow \exists v \, T(v, y) \wedge T'(v, z)$$

Σ

Schema  Source S   Target T

Database   I ⇢ J

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

# A Data Exchange Example

[Popa et al. 2001]

**S**: *Rcd*
  **projects**: *Set of*
    **project**: *Rcd*
      **name**
      **status**

  **grants**: *Set of*
    **grant**: *Rcd*
      **gid**
      **project**
      **recipient**
      **manager**
      **supervisor**

  **contacts**: *Set of*
    **contact**: *Rcd*
      **cid**
      **email**
      **phone**

  **companies**: *Set of*
    **company**: *Rcd*
      **name**
      **official**

**T**: *Rcd*
  **projects**: *Set of*
    **project**: *Rcd*
      **code**
      **funds**: *Set of*
        **fund**: *Rcd*
          **fid**
          **finId**

  **finances**: *Set of*
    **finance**: *Rcd*
      **finId**
      **mPhone**
      **company**

  **companies**: *Set of*
    **company**: *Rcd*
      **coid**
      **name**

**Projects**

| code: **E-services** |
| --- |

**Funds**

| fid | finId |
| --- | --- |
| g3 | **???** |

| code: **PIX** |
| --- |

**Funds**

| fid | finId |
| --- | --- |
| g1 | **???** |
| g2 | **???** |

**Finances**

| finId | mPhone | company |
| --- | --- | --- |
| **???** | 3608679 | ??? |
| **???** | 3608776 | ??? |
| **???** | 3608600 | ??? |

**Companies**

| coid | name |
| --- | --- |
| ??? | AT&T |
| ??? | Lucent |

**Target instance**

**project**(na,st), **grant**(gid,na,re,ma,su), **contact**(ma,em,ph) →

**project**(na,FUND), **fund**(gid,finId), **finance** (finId,ph,company),

**company**(company, name)

# Mapping generation

$E_\mathbf{S}$:

**Company**

| IdCompany | Name | Town |
|---|---|---|
| 'C1' | 'AA' | 'Paris' |
| 'C2' | 'Ev' | 'Lyon' |

**Flight**

| Departure | Arrival | IdCompany |
|---|---|---|
| 'Lyon' | 'Paris' | 'C1' |
| 'Paris' | 'Lyon' | 'C2' |

**Travel Agency**

| IdAgency | Name | Town |
|---|---|---|
| 'A1' | 'TC' | 'L.A.' |

$E_\mathbf{T}$:

**Firm**

| Id | Name | Town |
|---|---|---|
| 'Id1' | 'AA' | 'Paris' |
| 'Id2' | 'Ev' | 'Lyon' |
| 'Id3' | 'TC' | 'L.A.' |

**Departure**

| Town | IdFirm |
|---|---|
| 'Lyon' | 'Id1' |
| 'Paris' | 'Id2' |

**Arrival**

| Town | IdFirm |
|---|---|
| 'Paris' | 'Id1' |
| 'Lyon' | 'Id2' |

$\mathbf{m} : Company(c1, aa, paris) \wedge Company(c2, ev, lyon) \wedge TravelAgency(a1, tc, la)$

$\wedge Flight(lyon, paris, c1) \wedge Flight(paris, lyon, c2)$

$\rightarrow \exists id1, id2, id3, Firm(id1, aa, paris) \wedge Departure(lyon, id1) \wedge Arrival(paris, id1)$

$\wedge Firm(id2, ev, lyon) \wedge Departure(paris, id2) \wedge Arrival(lyon, id2) \wedge Firm(id3, tc, la)$

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

# Interactive Mapping

**Input**: set of examples
$$(E_S^1, E_T^1) \dots (E_S^n, E_T^n)$$

Normalization

$\Sigma_{norm}$

Question

Yes/No

Atom refinement

$\Sigma_{atRef}$

Question

Yes/No

Join refinement

$\Sigma_{final}$

**Output**: refined mapping

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Atom Refinement

**Input**: set of examples
$(E_S^1, E_T^1) \ldots (E_S^n, E_T^n)$

Normalization

$\Sigma_{norm}$

Question ← Atom refinement

Yes/No →

$\Sigma_{atRef}$

Ask the user and refine the left part of the rule

$$\{C_1; C_2; F_1; F_2; TA\}$$

$$\{C_1; C_2; F_1; TA\} \quad \{C_1; C_2; F_1; F_2\} \quad \{C_1; F_1; F_2; TA\}$$

$$\{C_1; C_2; TA\} \quad \{C_1; C_2; F_1\} \quad \{C_1; F_1; TA\} \quad \{C_1; C_2; F_2\} \quad \{C_1; F_1; F_2\} \quad \{C_1; F_2; TA\}$$

$$\boxed{\{C_1; C_2\}} \quad \{C_1; F_1\} \quad \{C_1; F_2\}$$

Are the tuples Company(c1,aa,paris); Company (c2, ev, lyon) enough to produce Firm(id, aa, Paris); Departure (Lyon, id); Arrival(Paris, id)?

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

# Atom Refinement

**Input**: set of examples
$(E_S^1, E_T^1) \ldots (E_S^n, E_T^n)$

Normalization
$\Sigma_{norm}$

Question

Yes/No

Atom refinement

$\Sigma_{atRef}$

Ask the user and refine the left part of the rule



$\{C_1; C_2; F_1; F_2; TA\}$

$\{C_1; C_2; F_1; TA\}$  $\{C_1; C_2; F_1; F_2\}$  $\{C_1; F_1; F_2; TA\}$

$\{C_1; C_2; TA\}$  $\{C_1; C_2; F_1\}$  $\{C_1; F_1; TA\}$  $\{C_1; C_2; F_2\}$  $\{C_1; F_1; F_2\}$  $\{C_1; F_2; TA\}$

$\{C_1; C_2\}$  $\{C_1; F_1\}$  $\{C_1; F_2\}$

Are the tuples Company(c1,aa,paris); Flight (lyon, paris, c1) enough to produce Firm(id, aa, Paris); Departure (Lyon, id); Arrival(Paris, id)?

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Where we are

Relational databases

Textual data

Graphs and networks

Machine learning

Challenges and Remarks

**SIGIR 2019** tutorial
https://data-exploration.ml

# SIMILARITY for DOCUMENTS

## Unstructured

## Semi-Structured

**SIGIR 2019** Tutorial

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

**SEARCHING** FOR

**Documents**

**Semi-Structured information**

BY **LOOKING AT**

| Words | Meta-Data | Words | Web-Tables |
|---|---|---|---|
| Text Classifier [Liu et al. '03 Zhang and Lee'09, Zhu et al.'13]<br><br>Topic Models [Zhu and Wu. '14]<br><br>Segmentation [Papadimitriou et al. '17] | Citation Graph Navigation [El-Arini et al.'11 Jia and Saule'17]<br><br>Entity Linking [Bordino et al. '13] | Regular Expressions [Agichtein et al.'00]<br><br>Annotations [Hanafi et al.'17]<br><br>Entity Extraction [Ritter et al.'15] | Entity Mentions [Wang et al.'15]<br><br>Schema Matching [Yakout et al.'18] |

**APPLYING**

**PRODUCES**

| Documents/Citations/Queries recommendations | Relation Extraction Document Matching | Entity Augmentation Concept Expansion |
|---|---|---|

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

# Document Search

Keyword Queries
& Relevance

Keyword query: search text with text
  **"Action movie with magic"**

Search documents containing those exact words
  … a live <u>action movie</u>…

  …. there is plenty of <u>action</u>…

  … packed <u>with action</u>…

  … <u>Magic</u> Mike is comedy <u>movie</u> …

  … in Harry Potter <u>magic</u> is everywhere..

Is this enough?
*Identify "relevant words"
and "relevant documents"*

**IMDb**

★★☆☆☆ **Super Mario Bros The Movie**
By Kay E. Platt on February 23, 2009
Hello People, I am going to be reviewing a Movie that ruined my school reputation…. The Movie itself is OK….

These famous actors who are chosen to play mario and luigi are acting in this movie, OK.. So I was in first grade when I watched this on VHS, and then my best friend Louis who was sitting next to me at story time was talking to me and then th…

★☆☆☆☆ **There are no magicians in this movie**
May 26, 2018
Format: DVD
I don't mean to give any spoilers away, but there are no magicians in this movie. Don't let the title fool you.

★☆☆☆☆
September 21, 2018
Format: Prime Video
Maybe don't name your musical "Rent" if you don't even have a single song about leasing law, property management procedures, or net lease calculations. As a real estate professional I am very disappointed and feel I was misled.

★☆☆☆☆ **Don't be gullible**
January 9, 2019
Format: Prime Video
This movie is dumb. Neil Armstrong was not very smart at all and Ryan playing him is just wrong. This guy (Armstrong) was not a successor at all. I believe that there are some critical information that doesn't quite add up to whether there was a …ecially since what more then … accomplish again. Why is …ite advanced today. I feel …e to reach something that is … and Mars when technology … I do not give Armstrong …y millions of Americans do.

★★★★★ **pokemon**
January 17, 2013
**Verified Purchase**
Format: VHS Tape
I will watch this while wearing pokemon clothes, sitting with my pokedoll, listening to the theme song, while playing pokemon on my ds.

# Document Search

## Relevant Keywords

**Relevance:** which keywords are more helpful in describing the content of the document?

**Relevance ≠ Frequency**

What keywords are more _likely to be used_ to describe the document we want and _not other documents_

1. Term-frequency: how many times the term appears in the document
2. Document-frequency: In how many documents the term appears

TF-IDF: Term Frequency Inverse Document Frequency



| MOANA | |
|---|---|
| Frequent | TF - IDF |
| film | maui |
| moana | te |
| the | moana |
| million | fiti |
| disney | cravalho |
| maui | goddess |

| THE INCREDIBLES | |
|---|---|
| Frequent | TF - IDF |
| film | parrs |
| the | syndrome |
| incredibles | violet |
| bird | omnidroid |
| pixar | parr |
| release | mirage |
|  | anisan |
|  | elen |

| MONSTERS INC. | |
|---|---|
| Frequent | TF - IDF |
| film | sulley |
| sulley | waternoose |
| monsters | boo |
| the | cda |
| mike | randall |
| monster | scarer |
| pixar | fizt |
| story | celia |

**FEW SELECTED KEYWORDS IN THE USER QUERY**
_What keywords to choose?_

**TRADITIONAL SEARCH**

**EXPLORATORY SEARCH**

**SEARCHING** FOR

**Documents**

**Semi-Structured information**

BY **LOOKING AT**

| Words | Meta-Data | Words | Web-Tables |
|---|---|---|---|

**APPLYING**

| | | | |
|---|---|---|---|
| Text Classifier [Liu et al. '03 Zhang and Lee'09, Zhu et al.'13] | Citation Graph Navigation [El-Arini et al.'11 Jia and Saule'17] | Regular Expressions [Agichtein et al.'00] | Entity Mentions [Wang et al.'15] |
| Topic Models [Zhu and Wu. '14] | | Annotations [Hanafi et al.'17] | |
| Segmentation [Papadimitriou et al. '17] | Entity Linking [Bordino et al. '13] | Entity Extraction [Ritter et al.'15] | Schema Matching [Yakout et al.'18] |

**PRODUCES**

| Documents/Citations/Queries recommendations | | Relation Extraction Document Matching | Entity Augmentation Concept Expansion |
|---|---|---|---|

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Documents as Examples

Exemplar documents

Set of exemplar documents

rather than a set of keywords.

**An entire document may contain more information!**
*It also contains more noise*

*Identify what makes them special, i.e., relevant*

**Example-based Document Search**
Given a corpus of documents D,
and a small set of relevant documents ($D_{rel}$),
identify a set of answer documents $D_A$
such that $D_{rel} \subseteq D_A \subseteq D$.

**Model as a classification problem**

**Find me movies like these:**



Monsters, Inc.
2001 · Fantasy/Adventure · 1h 32m

Monsters Incorporated is the largest scare factory in the monster world, and James P. Sullivan (John Goodman) is one of its top scarers. Sullivan is a huge, intimidating monster with blue fur, large purple spots and horns. His scare assistant, best friend and roommate is Mike Wazowski (Billy Crystal), a green, opinionated, feisty little one-eyed monster. Visiting from the human world is Boo (Mary Gibbs), a tiny girl who goes where no human has ever gone before.

The Incredibles
2004 · Action/Adventure · 1h 56m

In this lauded Pixar animated film, married superheroes Mr. Incredible (Craig T. Nelson) and Elastigirl (Holly Hunter) are forced to assume mundane lives as Bob and Helen Parr after all super-powered activities have been banned by the government. While Mr. Incredible loves his wife and kids, he longs to return to a life of adventure, and he gets a chance when summoned to an island to battle an out-of-control robot. Soon, Mr. Incredible is in trouble, and it's up to his family to save him.

**PROBLEM:MISSING NEGATIVE CLASS**
***Few positive examples*** and
a ***large set of unknown.***
What *features* can *discriminate* relevant and irrelevant?
Would be better to have *some negative examples*

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

# Text Classifiers

Using Positive and Unlabeled Examples

**Positive Unlabeled learning**

- a corpus of documents D,

- 2 Classes: relevant ⊤ & irrelevant ⊥

- relevant documents ($D_{rel}$)
  $\forall d \in D_{rel}$. class(d) = ⊤

- Unlabeled documents U = D − $D_{rel}$

**Goal:**
> train a classifier C : D → {⊤, ⊥},
> to predict class(u) ∀ u ∈ U.

**Missing:**
To train C we need examples
for the negative class ⊥

**Algorithm 4.9** Document Classification with Positive and Unalabled Data

**Input:** Relevant Documents $\mathbf{D}_{rel} \subseteq \mathcal{D}$, Unlabeled Documents $\mathbf{U} \subseteq \mathcal{D}$
**Output:** Classifier $\mathbb{C}$
1: $\mathbf{D}_{neg} \leftarrow$ getNegativeSample($\mathbf{U}$) ▷ See Li and Liu [2003], Liu et al. [2002], Yu et al. [2002]
2: $\mathbb{C} \leftarrow$ trainClassifier($\mathbf{D}_{rel}, \mathbf{D}_{neg}, \mathbf{U} \setminus \mathbf{D}_{neg}$) ▷ E.g., Expectation Maximization, SVM, or Rocchio
3: **return** $\mathbb{C}$

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

# Inferring Negative Examples (I)

Assign a label to Unlabeled data:

how to determine a negative sample set without asking the user

**4 Alternative approaches**

- **Naïve Bayes** (McCallum et al. [1998])
  - All unlabeled data are assumed negatives
  - NB-Classifier estimates **P(c|d)** based on on **P(w|c)** with $c \in \{\top, \bot\}$, $d \in D$, and words $w \in W$

- The **Rocchio** technique (Raskutti et al. [2002])
  - $\forall d \in D$  $\vec{d}$ is the TF-IDF vector representation
  - Build prototype vectors $\vec{c}_\top$ for documents in $D_{rel}$
  - and $\vec{c}_\bot$ for documents in U
  - Compare each $\forall d \in U$ with $\vec{c}_\top$ and $\vec{c}_\bot$
  - assign the class of the most similar vector

**Goal:**
Determine set of elements to be regarded as reliable negatives (RN)

**Train a "simplistic" classifier**

$$\vec{c}_\top = \alpha \frac{1}{|\mathbf{D}_{rel}|} \sum_{d \in \mathbf{D}_{rel}} \frac{\vec{d}}{\|\vec{d}\|} - \beta \frac{1}{|\mathbf{U}|} \sum_{d \in \mathbf{U}} \frac{\vec{d}}{\|\vec{d}\|}$$

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

# Inferring Negative Examples (II)

Assign a label to Unlabeled data:

how to determine a negative sample set without asking the user

**Goal:**
Determine set of elements to be regarded as reliable negatives (RN)

**Train a "simplistic" classifier**

**4 Alternative approaches**

- The **Spy** technique (Liu et al. [2002])
  - Extract a sample S from the positive example
  - Merge S in U (deploy the spies!)
  - Build NB classifier with EM
  - Determine threshold t such that all spies are correctly classified
  - Document above the threshold are considered negative

- **1-DNF**\* technique (Yu et al. [2002]).
  - *Disjunctive Normal Form*
  - *Positive Example Based Learning*
  - Get words $W_f \subset W$. $freq(w, D_{rel})/|D_{rel}| > freq(w, U)/|U|$
  - Remove from U all documents containing any word in $W_f$



$D_{rel}$

U

RN

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Training the Expert Classifier

Exploit the partial-supervision

**Expert Classifier**
Builds on the result of the first step to train a much more sophisticated and precise classifier.

- **1-shot approach**
  - Use $D_{rel}$ and RN and train a classifier (SVM or EM)

- **Iterative approach**
  - Use $D_{rel}$ and RN and train a classifier $C_i$
  - Use $C_i$ and extract new negative documents Q
  - Add Q to RN, train a new classifier $C_{i+1}$
  - Continue until no more negative documents are retrieved

  [*Optionally*] evaluate the last trained classifier over $D_{rel}$ and discard it if it performs poorly

Methods perform **poorly** when **the initial set of documents is very small**

**The Rocchio approach + EM is best for this case**

Advanced models with TF-IDF or Topic models
*Zhu et al. [2013] - Zhu and Wu [2014]*

Beware of Class Imbalance!
*SMOTE: Synthetic Minority Over-sampling Technique*

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

# Document Segmentation

Intention-based relatedness

**Model documents as Composite Objects**
Do not perform matching across the posts as a whole but across
*_fragments_* of them that are *_written for the same intention_*

**Intuition:**
Different parts of the document
Have different Purposes:
- Provide background information
- Describe Problem
- Ask question…

Extra RAID disk drives seem to be the solution to my problem but does adding RAID drives requires a **reformat and rebuild of the system to improve performance?**

**Doc C**

My boss gave me yesterday an HP Pavilion computer with Intel Matrix Storage System, a 320GB drive and Linux pre-installed. **I am thinking to add an extra disk drive using a RAID 0 or 1. Can I do it without having to rebuild the entire system?** I have already looked at the HP official web site for how to use a JBOD. But I have not found anything related to it.

**Doc B**

I have an HP system with a RAID 0 controller and 4 disks in form of a JBOD. I would like to install Hadoop with a replication 4 HDFS and only 320GB of disk space used from every disc. **Do you know whether it would perform ok or whether the partial use of the disk would degrade performance.** Friends have downloaded the Cloudera distribution but it didn't work. It stopped since the web site was suggesting to have 1TB disks. I am asking because I do not want to install Linux and then realize that my **hardware configuration is not the right one.**

**Doc A**

My HP Pavilion stops working after 15 min of activity. I called our technical department but no luck. Despite the many calls, I did not manage to find **a person with adequate knowledge to find out what is wrong.** All they said is bring it to up and we will see, which frustrated me. At the end I had the brilliant idea to move it to a cooler place and voila. No more p

**Doc D**

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Segmentation

Boundaries

**Communication means & Text Features**

Use text characteristics and identify points in which a significant variation of these characteristics occurs, and place a segmentation border there.

| | | | |
|---|---|---|---|
| Tense($CM_{tense}$) | present | past | future |
| Subject ($CM_{subj}$) | I/we | you | it/they/(s)he |
| Style ($CM_{qneg}$) | interrog. | negative | affirmative |
| Status ($CM_{pasact}$) | passive | active | |
| Part of Speech($CM_{pos}$) | verb | noun | adj./adverb |

0 **I** have an HP system with a RAID 0 controller and 4 disks in form of a JBOD. 75 **I would like** to install Hadoop with a replication 4 HDFS and only 320GB of disk space used from every disc. 182 Do **you know** whether 201 **it** would perform ok or whether the partial use of the disk 259 would degrade performance. 285 **Friends** have downloaded the Cloudera distribution but 338 **it** didn't work. 355 **It** stopped since 371 **the web site** was suggesting to have 1TB disks. 418 **I am asking** because 436 **I** do not want to install Linux and then realize that 488 **my hardware configuration** is not the right one. 535

**Good border**

**Not good border**

**segment 2**

$b_1$    $b_2$    $b_3$

| S $_{i-2}$ | S $_{i-1}$ | S $_i$ | S $_{i+1}$ |
|---|---|---|---|

**segment 1**

Bottom-up approach
1. Start with single words as segments
2. Compute a **Diversity Index** in each segment
3. Merge segments with low diversity

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Intention Clustering & Matching

~~Match...~~ ....ents with the same intention

**Clusters are based on intention**

Given a document $d_q$,
1. the system will **segment** $d_q$,
2. identify for each segment the **segments in the same cluster**
3. **aggregate the similarity** of those segments into a score for each document.

**C3**

Friends have downloaded the Cloudera distribution but it didn't work. It stopped since the web site was suggesting to have 1TB disks.

I am thinking to add an extra disk drive using a RAID 0 or 1. Can I do it without having to rebuild the entire system?

Do you know whetherit would perform ok or whether the

Despite the many calls, I did not manage to find a person with adequate knowledge to find out what is wrong.

**Explore based on related topics linked to common goals**

**C1**

I am asking because I do not want to install Linux and then ...lize that I need a hardware configuration i.e., the right one.

I have already looked at the HP official web site for how to use a JBOD. But I have not found anything related to it.

I have an HP system with a RAID 0 controller and 4 disks in

Extra RAID disk drives seem to be the solution to my problem but does adding RAID drives requires a reformat and rebuild of the system to improve performance?

**C2**

All they said is bring it to up and we will see, which frustrated me. At the end I had the brilliant idea to move it to a cooler place and voila. No more problems.

My HP Pavilion stops working after 15 min of activity. I called our technical department but no luck.

Linux pre-installed.

https://data-exploration.ml

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

**Document Networks**

Fast Algorithms for Mining Association Rules in Large Databases

Management of probabilistic data: foundations and challenges

Privacy-Preserving Data Publishing

Efficient Query Evaluation on Probabilistic Databases

Privacy-preserving data mining

The boundary between privacy and utility in data publishing

Privacy preserving OLAP

Rakesh Agrawal

Nilesh Dalvi

Dan Suciu

Sungho Hong

Vibhor Rastogi

Ramakrishnan Srikant

Dilys Thomas

**SIGIR 2019** tutorial          M. Lissandrini, D. Mottin, T. Palpanas, Y. Velegrakis

https://www.microsoft.com/en-us/research/project/microsoft-academic-graph/

# Influence in Citation Networks

El-Arini and Guestrin [2011]
Jia and Saule [2017]

Document relevance based on influence

## Citation Network

- Nodes are Authors and Papers

- Edges are Authorship and Citations

- Influence is based on connecting Paths

## Advance Models

- El-Arini and Guestrin [2011] :

  - Condition influence on topics
    *Iterate for each topic T: Select topic T, keep only papers relevant for T, compute connecting Paths.*

  - Weight edges with Influence-Probability

- Jia and Saule [2017]

  - Enrich graph with Keywords & Venues

Rakesh Agrawal

Privacy-preserving data mining

Ramakrishnan Srikant

The boundary between privacy and utility in data publishing

Dan Suciu

*Just looking at citations and co-citations is not sufficient.*

**Start from a known document Explore new related topics, authors, venues…**

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Traverse (Document) Networks

El-Arini and Guestrin [2011]
Jia and Saule [2017]



## Personalized Page Rank

- Start from seed nodes, i.e. the documents $D_{rel}$

- Navigate towards locally connected nodes

### Example based Exploration implies locality

### Global Page Rank

Starting from a random node, traversing randomly, **random restart point** anywhere in the graph

### CHALLENGE:
*Identify meaningful transition probabilities*

*E.g., El-Arini and Guestrin [2011]*

### Personalized Page Rank

Starting from a **limited set of nodes**, traversing randomly, restart point is one in **the initial set**. <u>Bound not to travel too far</u>

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

# Serendipitous Search

Enhance document links with Entities and Query-logs

**Input:** Query/Document
**Output:** Queries

America    Peru

Machu Picchu

Connected entities

Query Logs

Francisco Pizarro
Rafting excursion
Amazon river
…

**Serendipity**
*Related topics* potentially come to mind *after* consulting the page.

### Document (Wikipedia article: Machu Picchu)

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**

**Exploit "lateral connections" in User Search Behaviors**

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

Searches related to Document content

# Entity Query Graph

Bordino et al. [2013]

Entity-Query graph from queries to entities and back

**EQGraph Weighted Edges**

Queries in the same session

1. **query to query:**

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

Frequency-based approach

2. **entity to query**

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

The more queries entities share the higher the probability

3. **entity to entity**

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

Based on query to query edges

Personalized PageRank to score suggested queries

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

**SEARCHING** FOR

**Documents**

**Semi-Structured information**

BY **LOOKING AT**

| Words | Meta-Data | Words | Web-Tables |
|---|---|---|---|

**APPLYING**

Text Classifier
[Liu et al. '03
Zhang and Lee'09,
Zhu et al.'13]

Topic Models
[Zhu and Wu. '14]

Segmentation
[Papadimitriou et al. '17]

Citation Graph Navigation
[El-Arini et al.'11
Jia and Saule'17]

Entity Linking
[Bordino et al. '13]

Regular Expressions
[Agichtein et al.'00]

Annotations
[Hanafi et al.'17]

Entity Extraction
[Ritter et al.'15]

Entity Mentions
[Wang et al.'15]

Schema Matching
[Yakout et al.'18]

**PRODUCES**

| Documents/Citations/Queries recommendations | Relation Extraction Document Matching | Entity Augmentation Concept Expansion |
|---|---|---|

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

# Entity Mentions &Web-Tables

Documents &
semi-structured information

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

| HR Information | | Contact | |
|---|---|---|---|
| Position | Salary | Office | Extn. |
| Accountant | $162,700 | Tokyo | 5407 |
| Chief Executive Officer (CEO) | $1,200,000 | | |
| Junior Technical Author | $86,000 | | |
| Software Engineer | $132,000 | | |
| Software Engineer | $206,850 | | |
| Integration Specialist | $270,000 | | |

| State | Capital | Population | Largest City | Population |
|---|---|---|---|---|
| Alaska | Juneau | 31,275 | Anchorage | 291,826 |
| Alabama | Montgomery | 205,764 | Birmingham | 212,237 |
| California | Sacramento | 466,488 | Los Angeles | 3,792,621 |
| Connecticut | Hartford | 124,775 | Bridgeport | 144,229 |
| Delaware | Dover | 36,047 | Wilmington | 70,851 |
| Florida | Tallahassee | 181,376 | Jacksonville | 821,784 |
| Illinois | Springfield | 116,250 | Chicago | 2,695,598 |
| Kansas | Topeka | 127,473 | Wichita | 382,368 |

In fact, the *Chinese* NORP market has the *three* CARDINAL most influential names of the retail and tech space – *Alibaba* GPE, *Baidu* ORG, and *Tencent* PERSON (collectively touted as *BAT* ORG ), and is betting big in the global *AI* GPE in retail industry space . The *three* CARDINAL giants which are claimed to have a cut-throat competition with the *U.S.* GPE (in terms of resources and capital) are positioning themselves to become the 'future *AI* PERSON platforms'. The trio is also expanding in other *Asian* NORP countries and investing heavily in the *U.S.* GPE based *AI* GPE startups to leverage the power of *AI* GPE . Backed by such powerful initiatives and presence of these conglomerates, the market in APAC AI is forecast to be the fastest-growing *one* CARDINAL , with an anticipated *CAGR* PERSON of *45%* PERCENT over *2018 - 2024* DATE .

To further elaborate on the geographical trends, *North America* LOC has procured *more than 50%* PERCENT of the global share in *2017* DATE and has been leading the regional landscape of *AI* GPE in the retail market. The *U.S.* GPE has a significant credit in the regional trends with *over 65%* PERCENT of investments (including M&As, private equity, and venture capital) in artificial intelligence technology. Additionally, the region is a huge hub for startups in tandem with the presence of tech titans, such as *Google* ORG , *IBM* ORG , and *Microsoft* ORG .

| Structure | Country | City | Height (metres) | Height (feet) | Year built | Coordinates |
|---|---|---|---|---|---|---|
| Burj Khalifa | United Arab Emirates | Dubai | 828.1 | 2,717 | 2010 | 25°11′50.0″N 55°16′26.6″E |
| Tokyo Skytree | Japan | Tokyo | 634 | 2,080 | 2011 | 35°42′36.5″N 139°48′39″E |
| KVLY-TV mast | United States | Blanchard, North Dakota | 628.8 | 2,063 | 1963 | 47°20′32″N 97°17′25″W |
| Abraj Al Bait Towers | Saudi Arabia | Mecca | 601 | 1,972 | 2011 | 21°25′08″N 39°49′35″E |
| Lotte World Tower | South Korea | Seoul | 555.7 | 1823 | 2017 | 37°30′45″N 127°6′10″E |
| One World Trade Center | United States | New York, NY | 541 | 1,776 | 2013 | 40°42′46.8″N 74°0′48.6″W |
| Large masts of INS Kattabomman | India | Tirunelveli | 471 | 1,545 | 2014 | 8°22′42.52″N 77°44′38.45″E ; 8°22′30.13″N 77°45′21.07″E |
| Lualualei VLF transmitter | United States | Lualualei, Hawaii | 458 | 1,503 | 1972 | 21°25′11.87″N 158°08′53.67″E ; 21°25′13.38″N 158°09′14.35″W |
| | | Kuala Lumpur | 452 | 1,482 | 1998 | 3°09′27.45″N 101°42′40.7″E; 3°09′29.45″N 101°42′43.4″E |
| | | New York | 425.5 | 1,396 | 2015 | |

| | Country |
|---|---|
| | Germany |
| Pay Talk — Francisco Chang | Mexico |
| Earn More — Roland Mendel | Austria |
| Island Trading — Helen Bennett | UK |

| Texas | Austin | 790,390 | Houston | 2,099,451 |
|---|---|---|---|---|
| Virginia | Richmond | 204,214 | Virginia Beach | 437,994 |
| Vermont | Montpelier | 7,855 | Burlington | 42,417 |
| Washington | Olympia | 46,478 | Seattle | 608,660 |
| Wisconsin | Madison | 233,209 | Milwaukee | 594,833 |

# Entity-relation tuples

Example-based extraction of Entity mentions and Relations

Brin [1998]
Agichtein and Gravano [2000]

**Search for Information WITHIN Docume**
*Explore new Entities*
*and new ways to express relations*

Works bests with Binary relation
Can work with multiple mentions:

Bob born in U.S.A. in 1978

**1. Example**

⟨ Google ; Menlo Park ⟩

**2. Match**

Google founded in Menlo Park...

**3. Extract Pattern**

... [X] founded in [Y] ...

**4. Extract New Mentions & Patterns**

Apple founded in Coupertino ...

Apple headquarters in Coupertino

Exemplar Tuples

Find Occurences of Exemplar Tuples

Generate New Exemplar Tuples

**Snowball**

Tag Entities

Augment Table

Generate Extaction Patterns

**Goal:** Enrich a list of Entity-relationships data

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Entity-relation tuples

Example-based extraction of Entity mentions and Relations

Brin [1998]
Agichtein and Gravano [2000]

*How to validate the new rules extracted automatically?*

1. **Compare extracted rules with known tuples:** confidence of R is based on how many known tuples extracts

2. **Compare extracted tuples with known rules:** confidence of T is based on how many known rules also extract T

**New extracted Rules and Tuples should not create contradictions**

*This approach has no "human in the loop"*

Exemplar Tuples → Find Occurences of Exemplar Tuples → Tag Entities → Generate Extaction Patterns → Augment Table

Generate New Exemplar Tuples

**Snowball**

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# IN MY DEFENSE
# I WAS LEFT UNSUPERVISED



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

# Entity-extraction by Example

Hanafi et al., [2017]

Learn extraction rules from example

Allow to match from text
both **Positive** and **Negative** examples

**Goal:** Supervised Extraction

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

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Matching Rules

From string tokens to "semantics"

Hanafi et al., [2017]

**Intuition:** Exploit a vocabulary of simple specialized patters with known semantics

Example:  | 5 percent up in  Dubai |

| 5 |

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

| percent |

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

...

Each rule has a "class" and a preference score

Each token may have different candidate "matching rules"

Token gap Regex  <  Literal Dictionary  < Pre-builts

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

0 ◁◁◁◁◁◁◁◁ 1

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

# Merging Rules

Reconcile multiple interpretations

Example: 5 percent

Tokens: 5     percent

Tree:

- P: Percentage = 1.0
- L: '5' = 0.4
- L: 'percent' = 0.4
- R: [A-Za-z]+ ≡ 0.2
- 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 %

Tokens: 6     %

- P: Percentage = 1.0
- L: '6' = 0.4
- R: [0-9]+ ≡ 0.2
- L: '%' = 0.4
- R: symbols = 0.2
- L: '%' = 0.4
- R: symbols = 0.2

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

Intersection: [5 percent, 6%]

- P: Percentage = 1.0
- L: {'5', '6'} = 0.4 —— L: {percent, %} = 0.4
- R: [0-9]+ = 0.2 —— L: {percent, %} = 0.4

*Consider also Negative Examples to prune candidates*

**SIGIR 2019** tutorial
https://data-exploration.ml

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

https://vimeo.com/208729128

# Web Tables

Semi-structured data on the web



https://en.wikipedia.org/wiki/Denmark#Regions

## Regions

The governing bodies of the regions are the regional councils, each with forty-one councillors elected for four-year terms. The councils headed by regional district chairmen (*regionsrådsformanden*), who are elected by the council.[79] The areas of responsibility for the re councils are the national health service, social services and regional development.[79][80] Unlike the counties they replaced, the region allowed to levy taxes and the health service is partly financed by a national health care contribution until 2018 (*sundhedsbidrag*), partl from both government and municipalities.[18] From 1 January 2019 this contribution will be abolished, as it is being replaced by higher tax instead.

The area and populations of the regions vary widely; for example, the Capital Region, which encompasses the with the exception of the subtracted province East Zealand but includes the Baltic Sea island of Bornholm, has than that of North Denmark Region, which covers the more sparsely populated area of northern Jutland. Unde densely populated municipalities, such as Copenhagen Municipality and Frederiksberg, had been given a stat making them first-level administrative divisions. These *sui generis* municipalities were incorporated into the ne reforms.

| Danish name | English name | Admin. centre | Largest city (populous) | Population (January 2017) | Total area (km²) |
|---|---|---|---|---|---|
| Hovedstaden | Capital Region of Denmark | Hillerød | Copenhagen | 1,807,404 | 2,568.29 |
| Midtjylland | Central Denmark Region | Viborg | Aarhus | 1,304,253 | 13,095.80 |
| Nordjylland | North Denmark Region | Aalborg | Aalborg | 587,335 | 7,907.09 |
| Sjælland | Region Zealand | Sorø | Roskilde | 832,553 | 7,268.75 |
| Syddanmark | Region of Southern Denmark | Vejle | Odense | 1,217,224 | 12,132.21 |

**Source:** Regional and municipal key figures

89

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

# Entity List Expansion

Augmentation: identify entities to complete the list

**Goal:** Given some seed entity mentions, retrieve more entities of the same type

1. Input: Incomplete list + Keyword query
2. Retrieve tables from pages based on the keyword query
3. Assign Score to tables based on relevance
4. Extract entity mentions from tables
5. Analyze Entity mention co-occurrence
6. Pick "co-occurring" Entities

Google
Accenture          0.4
Facebook
Dell               0.3
IBM
HP                 0.5
**Bipartite-graph**
                   0.6

**Problem:** *entities may appear together for different reasons*

**Score Propagation**

**Problem:** *Here PPR Causes concept drift*

**Heuristic Propagatio**

*Incomplete table*

| IT Company |
| --- |
| Dell |
| IBM |
| Lenovo |
| ….? |

*Augmented table*

| IT Company |
| --- |
| Dell |
| IBM |
| Lenovo |
| **Apple** |
| **Samsung** |
| **HP** |
| **Acer** |

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Web-Table Completion

Identify relevant content, retrieve missing information

**Goal:** Retrieve missing attribute values

**Intuition:** If there is a structure, we can match it!

| Model | Brand |
|---|---|
| S80 | N |
| Easyshare CD44 | K |
| DSC W570 | S |
| Optio E60 | P |

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

| Part N |
|---|
| DSC W |
| T146 |
| Optio E |
| S8100 | Nikon |

**Web tables**

**Incomplete table**

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

**InfoGather**

**Problem:** *entities may appear together for different reasons*

**Complete table**

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

**Extra Input: table header**
target attribute name or example of completing attribute

# Table Correlation Graph

Schema matching for web-page and web-tables
Binary-relations only

**Goal:** Retrieve missing attribute values

## Determine Table Match

Direct Match between Q(K,A) and T(K,B)
K=entity names in a column
A,B = attribute column name (header)

Can use approximate
matching and thesaurus

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

**Problem:** *considers only direct links between Q and T*

**Query**

| Name | Developer |
|------|-----------|
| SQL Server | Microsoft |
| MySQL | |
| Teradata | |
| Firebird | |

**T1**

| Product | Vendor |
|---------|--------|
| MySQL | Oracle corp. |
| PostgreSQL | PostgreSQL Grp |
| MongoDB | MongoDB inc. |
| Berkley DB | Oracle corp. |

List of Open Source **database software**

**T2**

| Name | Max Row Size |
|------|--------------|
| MySQL | 64Kb |
| Oracle | 8Kb |
| Firebird | 64Kb |
| Berkley DB | 8kb |

Information about **database** size limits

**T3**

| Name | Developer |
|------|-----------|
| MySQL | Oracle |
| SQL Server | Microsoft |
| Office | Micrsoft |
| Photoshop | Adobe |

Best selling **software** in 2010

**T6**

| Vendor | Software |
|--------|----------|
| Oracle corp. | Oracle DB |
| IBM | DB2 |
| Teradata | Teradata Corp. |

Companies developing **database software**

**T5**

| Vendor | Revenue |
|--------|---------|
| Oracle | 11787M |
| IBM | 4870M |
| Microsoft | 4098M |
| Teradata | 882M |

**Database software**, 2011 revenue by vendor

**T4**

| Name | Windows | Linux |
|------|---------|-------|
| Oracle | Yes | Yes |
| MySQL | Yes | Yes |
| SQL Server | Yes | No |
| PostgreSQL | Yes | Yes |

OS support for top **database software**

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

# Table Correlation Graph

Schema matching for web-page and web-tables
Binary-relations only

**Goal:** Retrieve missing attribute values

**Determine Table Match**

Holistic Match
1. Assign Direct Match Score from Query to Tables
2. Scores >0 are starting nodes

3. Use classifier to add weight to other table pairs

   **Build Classifier** using
   - Context similarity
   - Table-to-content similarity
   - URL similarity
   - Tuples Similarity
   the model predicts the match between
   two tables with a probability

4. Use starting node and execute PPR

5. Use PPR scores to rank matching tables

Overcomes problems due to poor matching with the query

**T1**

| Product | Vendor |
|---------|--------|
| MySQL | Oracle corp. |
| PostgreSQL | PostgreSQL Grp |
| MongoDB | MongoDB inc. |
| Berkley DB | Oracle corp. |

List of Open Source **database software**

**T6**

| Vendor | Software |
|--------|----------|
| Oracle corp. | Oracle DB |
| IBM | DB2 |
| Teradata | Teradata Corp. |

Companies developing **database software**

0.1

**Query**

| Name | Developer |
|------|-----------|
| SQL Server | Microsoft |
| MySQL | |
| Teradata | |
| Firebird | |

0.2

**T2**

| Name | Max Row Size |
|------|--------------|
| MySQL | 64Kb |
| Oracle | 8Kb |
| Firebird | 64Kb |
| Berkley DB | 8kb |

Information about **database** size limits

0.3

0.7

**T5**

| Vendor | Revenue |
|--------|---------|
| Oracle | 11787M |
| IBM | 4870M |
| Microsoft | 4098M |
| Teradata | 882M |

**Database software,** 2011 revenue by vendor

0.2

0.1

0.2

**T3**

| Name | Developer |
|------|-----------|
| MySQL | Oracle |
| SQL Server | Microsoft |
| Office | Micrsoft |
| Photoshop | Adobe |

Best selling **software** in 2010

0.1

0.5

**T4**

| Name | Windows | Linux |
|------|---------|-------|
| Oracle | Yes | Yes |
| MySQL | Yes | Yes |
| SQL Server | Yes | No |
| PostgreSQL | Yes | Yes |

OS support for top **database software**

**SEARCHING** FOR

**Documents**

**Semi-Structured information**

BY **LOOKING AT**

| Words | Meta-Data | Words | Web-Tables |
|---|---|---|---|
| Text Classifier [Liu et al. '03 Zhang and Lee'09, Zhu et al.'13]<br><br>Topic Models [Zhu and Wu. '14]<br><br>Segmentation [Papadimitriou et al. '17] | Citation Graph Navigation [El-Arini et al.'11 Jia and Saule'17]<br><br><br>Entity Linking [Bordino et al. '13] | Regular Expressions [Agichtein et al.'00]<br><br>Annotations [Hanafi et al.'17]<br><br>Entity Extraction [Ritter et al.'15] | Entity Mentions [Wang et al.'15]<br><br><br>Schema Matching [Yakout et al.'18] |
| Documents/Citations/Queries recommendations | | Relation Extraction Document Matching | Entity Augmentation Concept Expansion |

**APPLYING**

**PRODUCES**

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Where we are

Relational databases

Textual data

**Graphs and networks**

Machine learning

Challenges and Remarks

**SIGIR 2019** tutorial
https://data-exploration.ml

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

**Protein Interaction Network**

**Road Network**

**Knowledge Graph**

**Social Network**

# Graphs *are* Everywhere

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

# Graphs
## Connected Data
*A graph is a graph is a graph*

| Fact Graph |
| --- |

| Ontology Tree |
| --- |

Edge-labelled Multigraphs

$G$: $\langle V, E, L, \ell \rangle$

Attributes:

$V/E$:<key,value>

Arnold Schwarzenegger

is A

Person

actedIN

is A

subClassOf

Terminator

| Release | 1984 |
| --- | --- |

Actor

**RDF** **(subject,predicate,object)**

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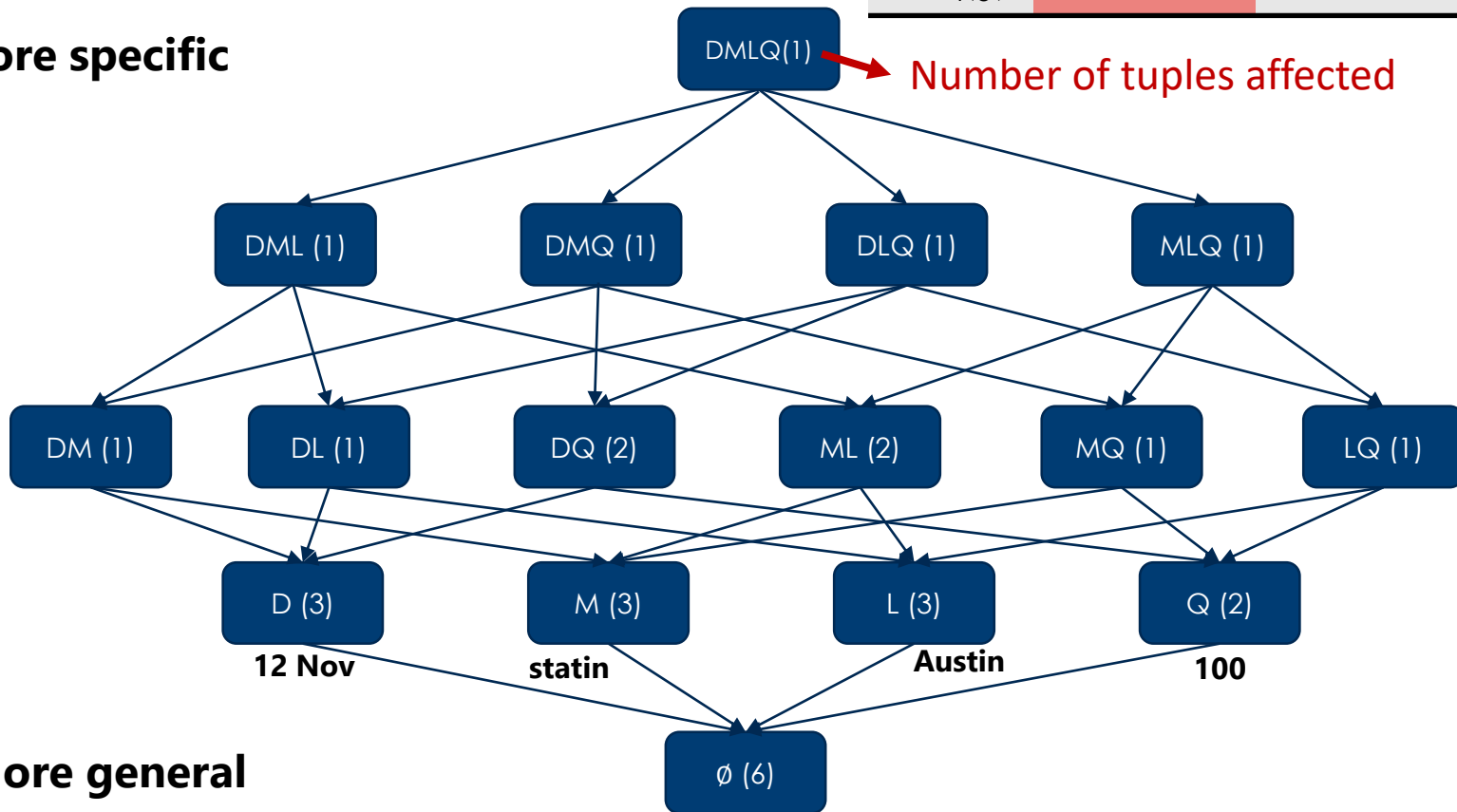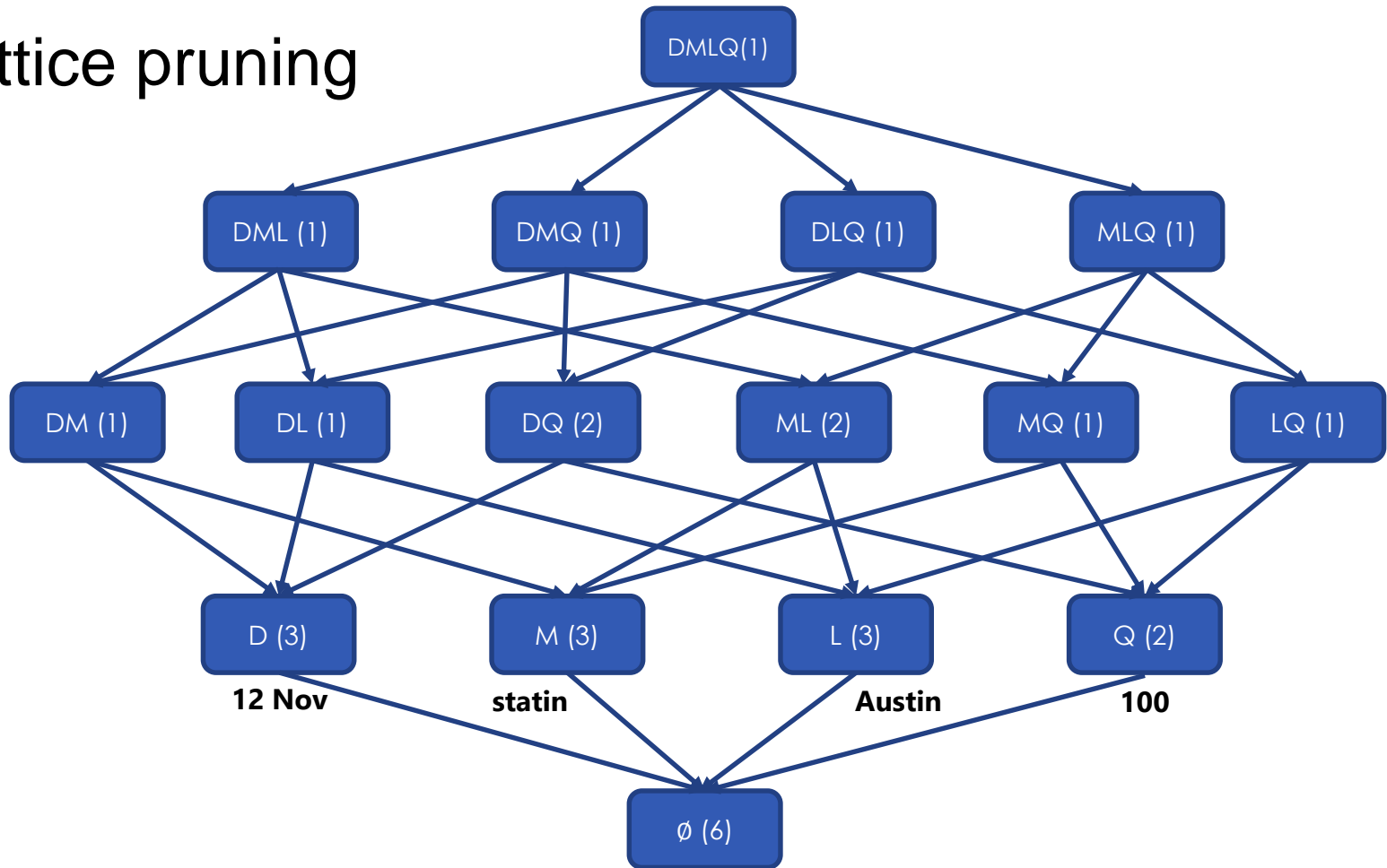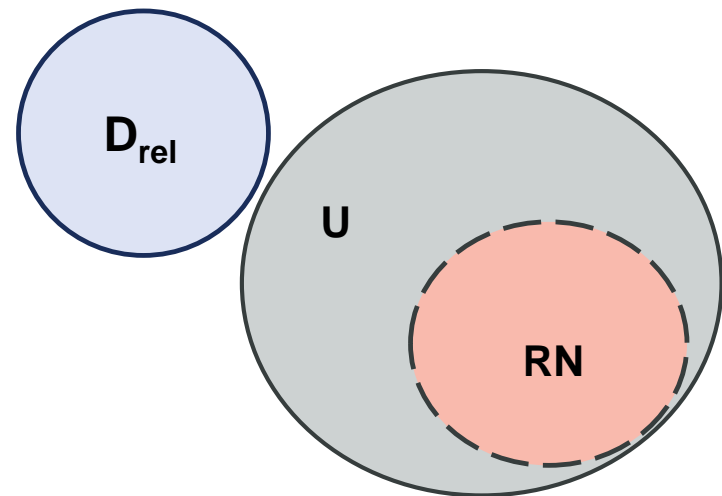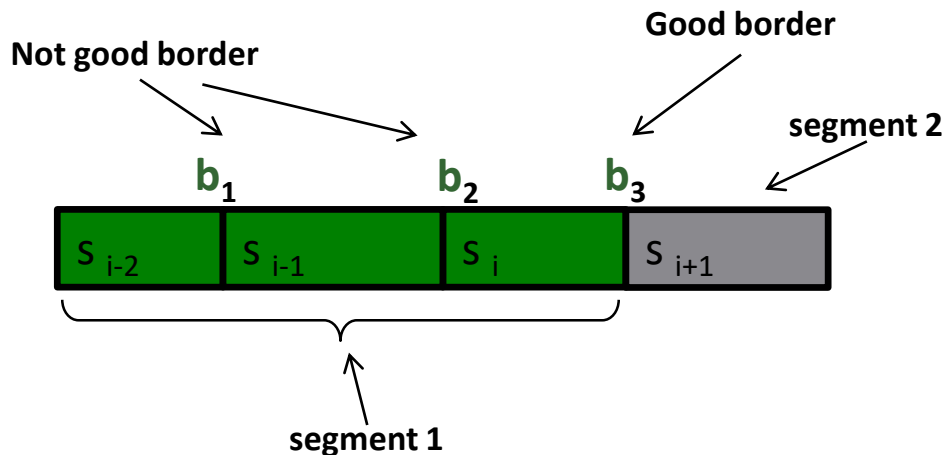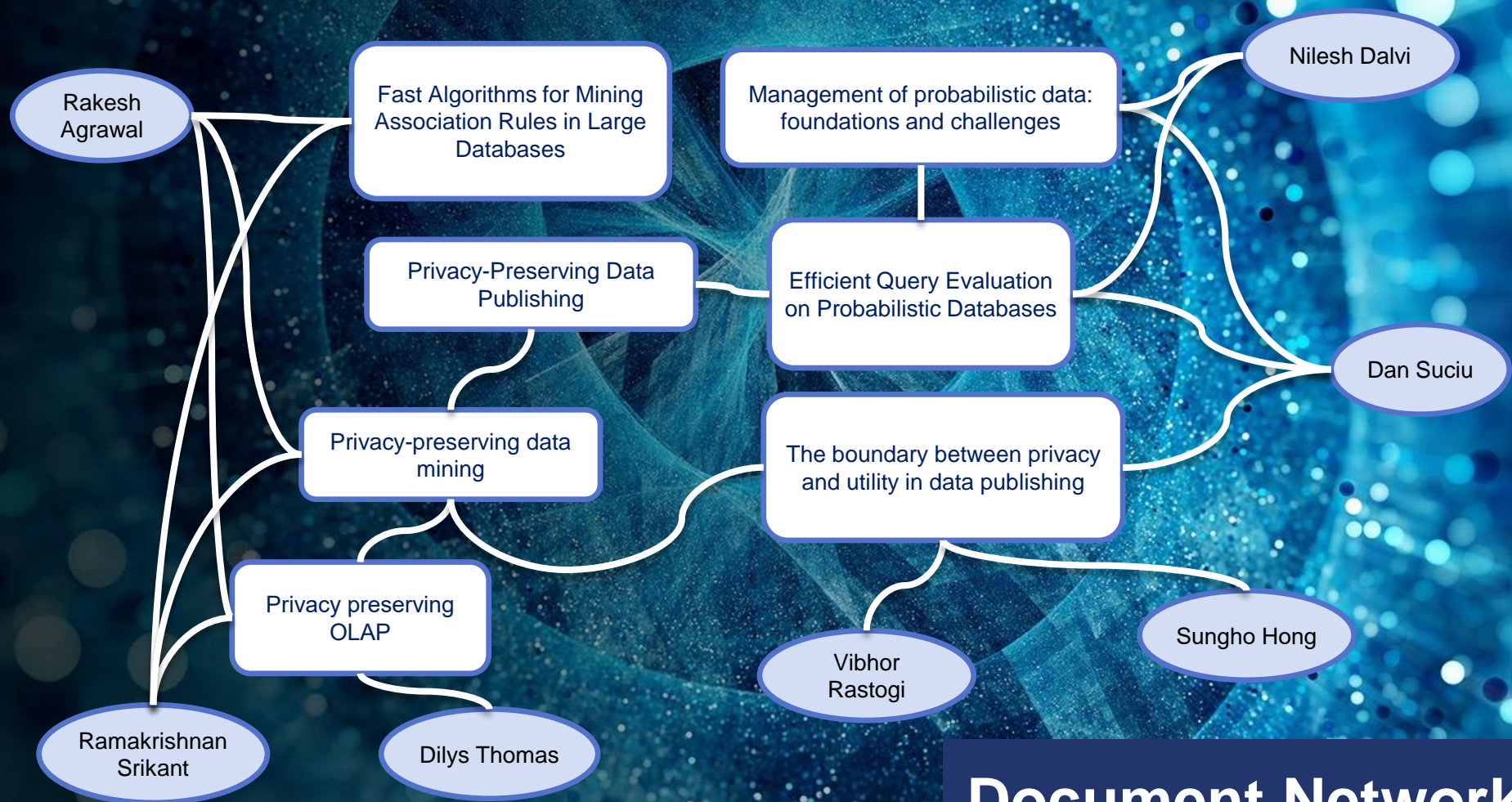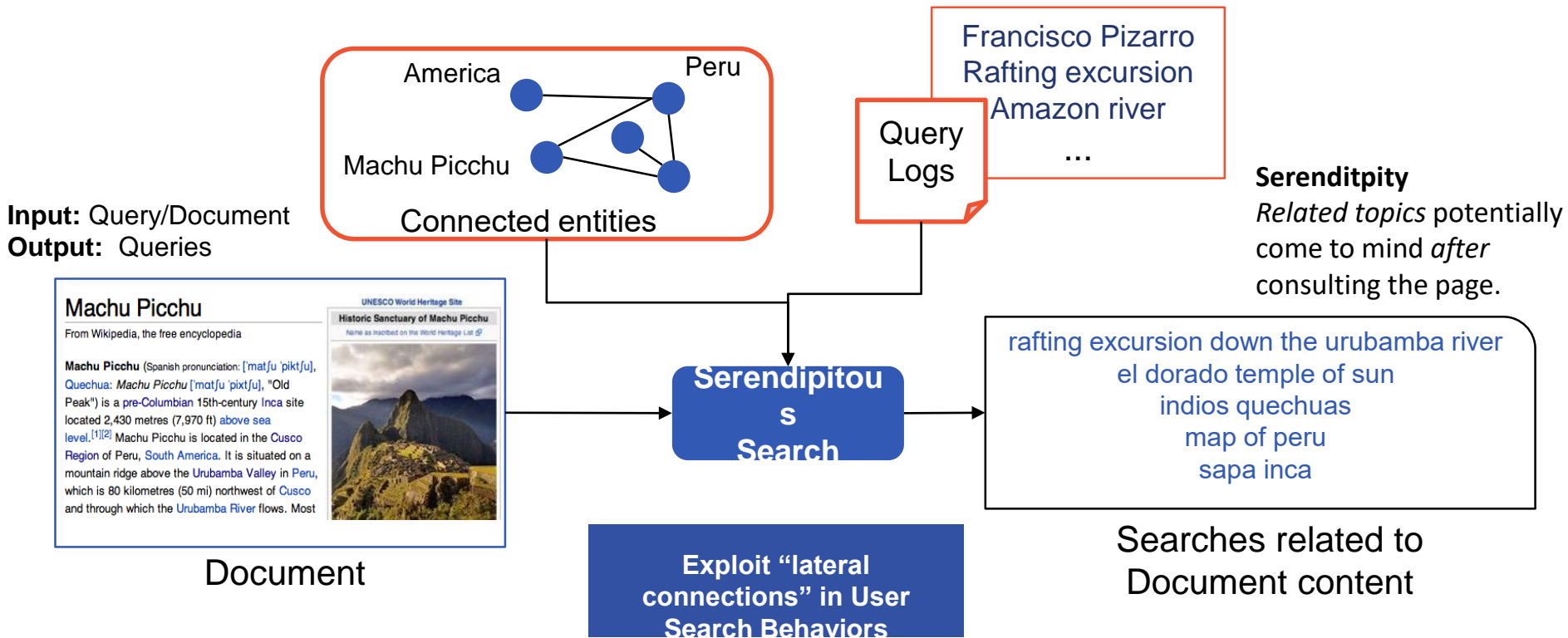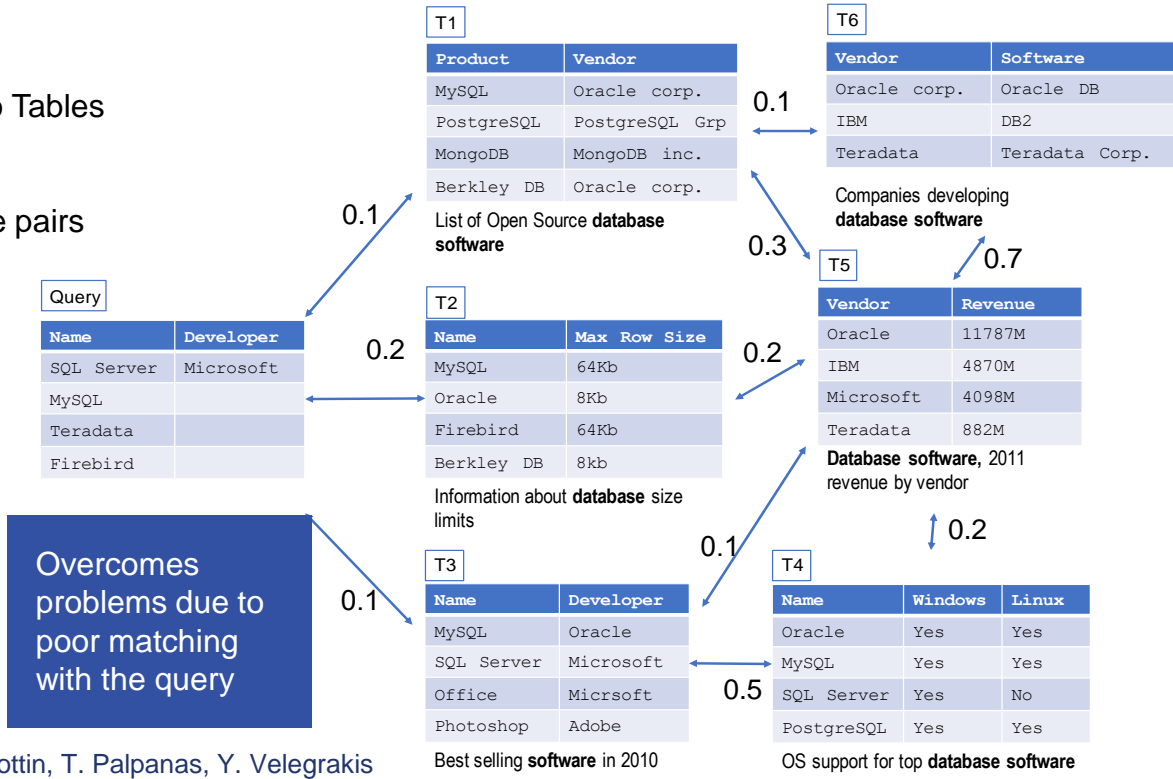(Arnold_Schwarzenegger,isA,Person)
      (Actor, subClassOf, Person)
(Arnold_Schwarzenegger, actedIn, Terminator)
```

**The Structure of the Graph is as important as the Data-values**

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

# Exemplar Queries

Example-driven graph search

**Input:** $Q_e$, an example <u>element</u> of interest

**Output:** set of <u>elements</u> 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 <u>relevant</u>

> Usually requires an intermediate step:
> User input (keywords) → Element in the graph

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

# SIMILARITY for GRAPHS

Nodes

Structures

| Connectivity | Properties |
| Queries | (Edge-)Labels |

**CHALLENGE: DISCOVER USER PREFERENCE**

**CHALLENGE: EFFICIENT SEARCH**

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

**SEARCHING** FOR

**Nodes**

**Structures**

BY **LOOKING** AT

Connectivity

Properties

Queries

(Edge-)Labels

**PRODUCES**

Mediator Nodes
[Gionis et al. '15
Ruchansky et al.'15]

Clusters
[Perozzi et al.'14,
Kloumann et al. 14]

Similar Entities
[Metzger et al.'13,
Sobczak et al.'15]

Path Queries
[Bonifati et al.'15]

SPARQL
[Arenas et al.'16]

Entity Tuples
[Jayaram et al.'15]

Similar
Structures
[Mottin et al.'14,
Xie et al.'17,
Lissandrini et al'18]

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

# Seed Set Expansion

Kloumann and Kleinberg [2014]

Nodes connected
by a community

Given a graph G, and a set of **query nodes** $V_Q \subseteq V_G$,
**retrieve all other nodes** $V_C \subseteq V_G$,
where C is a community in G, and $V_Q \subseteq V_C$.



Solution: PPR

$$\mathbf{v}^{t+1} = (1-\alpha)\mathbf{M} \cdot \mathbf{v}^t + \alpha \mathbf{v}^0$$

Communities can be <u>extremely large</u>
Identify "central nodes"
or "the core subgraph"

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

# The Minimum Wiener Connector Problem

**Model:**     Unlabeled Undirected Graph

**Query:**      A set of Nodes Q

**Similarity:**   Shortest-Path distance

**Output:**      A Set of <u>Connector Nodes</u> H

           "explains" connections in Q

Connectors:
Nodes with <u>HIGH</u> closeness
to ALL the inputs

Similar to a Steiner-Tree but
<u>overall pairwise distances</u> are optimized



Case: Infected Patients
→ Culprit/Other Infected

Case: Target Audience
→  Influencers

M. Lissandrini, D. Mottin, 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 <u>Connector Nodes</u> H

"explains" connections in Q

<span style="color:green">minimize the sum of pairwise shortest-path-distances between nodes in the connector H</span>

**Called: Wiener Index.**

tradeoff between size and average distance

W=1+2+1 =4

**W=1+1+1 = 3**

Sometimes The Best Solution is NOT A Tree

NP-Hard

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

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

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

# Approximate minimum Wiener Index Connector

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

All Pairwise Distances

→ Distances from a root r

Measure distance in H (i.e., subgraph-induced)

→ Precomputed distance in G

Edge Weights

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

**Approximated with Edge-Weighted SteinerTree**

**Enumerate Candidate Solutions for r ∈ Q & λ and keep best tree**



r

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Focused Clustering and Outlier Detection

## Similarity based on attributes

**Model:**      Unlabeled Undirected Graph <u>with Node Attributes</u>

**Query:**       A set of Nodes Q

**Similarity:** <u>To Be Inferred</u>

*based on Attribute Values & Connectivity*

**Output:**     <u>Clusters</u> of Nodes: Dense & Coherent

+ <u>Outliers</u>

Case: Target Users → Community with same interests

Case: Products→ Co-purchased products with similar features

Perozzi et al. [2014]

**PhD**
**NYC**
English
Google

~~College Paris Dutch Google~~

**PhD**
**NYC**
Greek
SAP

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

**PhD**
**NYC**
Italian
IBM

**PhD**
**NYC**
French
SAP

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Focused Clustering

## Infer User Focus

TASK: Infer "FOCUS" , important attributes

attribute weights β

$$\begin{pmatrix} \textbf{PhD} \\ \textbf{NYC} \\ \text{English} \\ \text{Google} \end{pmatrix} \quad \begin{pmatrix} \textbf{PhD} \\ \textbf{NYC} \\ \text{French} \\ \text{SAP} \end{pmatrix} \Rightarrow \begin{pmatrix} \textbf{0.5} \\ \textbf{0.5} \\ 0 \\ 0 \end{pmatrix}$$

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

$$\min_A \sum_{(u,v) \in P_S} (f_i - f_j)^T \mathbf{A} (f_i - f_j) - \gamma \log \left( \sum_{(u,v) \in P_D} \sqrt{(f_i - f_j)^T \mathbf{A} (f_i - f_j)} \right)$$

( Distance Metric Learning,   inverse Mahalanobis distance: Xing, et al 2002)



**Weak Ties**

Similar Attributes
=
Stronger Connection

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

# Focused Clustering

Perozzi et al. [2014]

Prune the Graph and keep dense communities

**TASK: Extract Clusters on Focused Graph**

attribute weights β → Edge Weight

**1. Find Starting Set** of Small Candidate Clusters

    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: $\varphi^{(w)}(C, G)$

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

    2.c Prune Underperforming nodes

**3. Detect Outliers:** High <u>unweighted</u> conductance

    w.r.t. low weighted conductance

LOCAL clusters

?

Seed

**Weighted Conductance:**
*ratio between the weighted sum of edges crossing the boundaries of the cluster and the weighted sum of those residing within it.*

**Performant Strategy:**
Start with local solution and expand around them to avoid complete scans of the graph

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

# iQBEES: Entity Search by Example

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

## Knowledge Graph Search

**Model:**      Knowledge Graph (Edge-labels)

**Query:**      A set of Entities Q

**Similarity:**    Shared semantic properties

**Output:**      A Set of <u>Similar Entities</u> (ranked)

Case: Products→ Products with similar aspects

Case: Social Media→ User recommendation

Entity 1:

Entity 2:

?

?

?



  **SIGIR 2019** tutorial
https://data-exploration.ml      M. Lissandrini, D. Mottin, T. Palpanas, Y. Velegrakis

# Maximal Aspects
## Selecting Features of Entity Similarity

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

?x sport BodyBuilding

?x type AmericanActor

Is not maximal if
Adding any aspect
→ E(A)={Arnold}

?x type AmericanActor

?x governorOf California

Include
Typical Types

**1. Prune generic aspects**

?x hasHeight 1.88m

?x type Entity

Use most
Specific Type

**2. Rank Set of aspects**

?x type AmericanActor

?x actedIn TheExpendables

?x type ActionActor

REPEATABLE
Update Q

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

# SIMILARITY for GRAPHS

Nodes

Structures

✓ Connectivity

✓ Properties

Queries

(Edge-)Labels

**Queries can retrieve
both Nodes and Structures**

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

**SEARCHING** FOR

**Nodes**

**Structures**

BY **LOOKING** AT

| Connectivity | Properties | Queries | (Edge-)Labels |
|---|---|---|---|

**PRODUCES**

Mediator Nodes
[Gionis et al. '15
Ruchansky et al.'15]

Clusters
[Perozzi et al.'14,
Kloumann et al. 14]

Similar Entities
[Metzger et al.'13,
Sobczak et al.'15]

Path Queries
[Bonifati et al.'15]

SPARQL
[Arenas et al.'16]

Entity Tuples
[Jayaram et al.'15]

Similar
Structures
[Mottin et al.'14,
Xie et al.'17,
Lissandrini et al'18]

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

# Learning Path Queries on Graphs

## Queries from Examples

**Model:** Edge Labeled Graph

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

Positive, Negative

**Similarity:** Common Path Query (RegExp)

$$q := \epsilon \mid a(a \in \Sigma) \mid q_1 + q_2 \mid q_1 \cdot q_2 \mid q^*$$

(bus|tram)*+ Cinema

**Output:** Set of Nodes satisfying paths for $Q^+$

but not paths for $Q^-$

Negative Examples to disambiguate intention

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

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Learnability of Path Queries

When is possible and How

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

Sometimes Positive & Negative Examples Cannot be reconciled!

**Consistency:**

1. Select Smallest Consistent Path
$$\forall v \in Q^+. \, paths_G(v) \nsubseteq paths_G(Q^-)$$

2. Loops cause infinite paths? Fix Maximal Length K
   When to use Kleene star * ?

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

3. Generalize SCP
   a) Construct Prefix Tree Acceptor
   b) Generalize into DFA with Merge

Can be INTERACTIVE! The system presents to the user nodes to label as Positive/Negative

Consistency Check: PSPACE-complete

Enumerate Paths Up to Fixed distance

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



PTA

DFA

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Reverse engineering SPARQL queries

Arenas et al. [2016]

## Knowledge Graph Search

**Model:** Knowledge Graph (Edge-labels)

**Query:** Set of Answers → Not Graphs but Tuples (of Nodes?)

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

**Output:** a SPARQL query / query results

Case: Open Data→ Query Unknown Schema

Case: Novice User → Avoid SPARQL

Mexico ·······?······· Spanish

Haiti

Jamaica ····?···· English

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

**MATCH** (?X, is_a, Country)
**OPT** (?X, has_language, ?Y)

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Complex SPARQL queries

A quick-peek to the complex pattern queries

```
SELECT * WHERE {
    ?deal a :Deal ;
          :employee ?employee ;
          :customer ?customer .
       ?employee :name ?employeeName ;
              :involvedAsEmployee ?matter .
       ?customer :name ?customerName ;
              :involvedAsCustomer ?matter .
}
```



Variables start with ?

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Reverse engineering SPARQL queries

Arenas et al. [2016]

## Challenges and Complexity

**Query:**      Set of Variable Mappings

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

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

OPT        OPT

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

Incomplete Mappings are
treated as OPTIONAL
Typical of RDF queries

Enumerate all possible
SPARQL queries satisfied
by the mappings

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

Build tree-shaped
SPARQL queries IMPLIED
by the mappings

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

# Reverse engineering SPARQL queries

## Challenges and Complexity
Query:     Set of Variable Mappings $\Omega$

3 Instantiations:
1. Only Positive Examples
2. Positive & Negative
3. Exact Result only

$$\Omega$$

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

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

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

$\{M4\}\ ?W$

$?X$      $?X$

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

$?W$      $?W$

$D$

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

Greedy: keep just enough to cover all variables

$(?X,t,P)\ ?X$

OPT    OPT

$(?X,e,?Y)\ ?Y$    $(?X,a,?Z)\ ?Z$

OPT

$(?Y,d,?W)(b4,d,?W)\ (?Y,d,d4)\ ?W$

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# SEARCHING FOR

**Nodes**

**Structures**

✓

## BY **LOOKING** AT

| Connectivity | Properties | Queries | (Edge-)Labels |

## PRODUCES

Mediator Nodes
[Gionis et al. '15
Ruchansky et al.'15]

Clusters
[Perozzi et al.'14,
Kloumann et al. 14]

Similar Entities
[Metzger et al.'13,
Sobczak et al.'15]

Path Queries
[Bonifati et al.'15]

SPARQL
[Arenas et al.'16]

Entity Tuples
[Jayaram et al.'15]

Similar
Structures
[Mottin et al.'14,
Xie et al.'17,
Lissandrini et al'18]

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

# Graph Exemplar Queries

## Search for Structures

**Model:** Knowledge Graph

**Query:** Example Structure

**Similarity:** Isomorphism/Simulation

**Output:** A set of Sub-Graphs



Case: Rich Schema ⟶ Find complex structures

# Graph Isomorphism vs. Simulation Variants

Structural Con...

Isomorphism requires an <u>bijective function</u>
Simulation requires only a <u>surjective relation</u>
Preserves only Parent → Child relationships



Example of **Simulating** (G1∼ {G2,G3,G4}) and **Strong-simulating** Graphs (G1≈G2)

**Strong Simulation preserves close connectivity**

Strong simulation: Capturing topology in graph pattern matching
– Shuai Ma et al., 2014

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

# Computing Exemplar Queries (i)

Mottin et al. [2016]

Fast Structure Matching

**Reduce Search Space:**
Removes nodes that **cannot be part of a** solution

NP-complete
(subgraph isomorphism)

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

**Exact Pruning technique:**
- Compute the neighbor labels of each node

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

- **Prune nodes not matching query** nodes **neighborhood** labels

- Apply iteratively on the query nodes

v A
B
Q

A
B
Sample

A
B
A1

A
B
A2

u

X

Labels at distance 1

neighborhood (v) = {(B,1)}

$\not\subseteq$

neighborhood (u) = {(A,1)}

No Match

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Computing Exemplar Queries (ii)

Mottin et al. [2016]

Prune Irrelevant Answers

**Reduce Search Space:**
Removes nodes that **are likely to be less relevant**

NP-complete
(subgraph isomorphism)

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

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

$$v = (1 - c)Av + cp$$

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

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

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



v

Sample   A1

A2

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

# Ranking Results
## Score Relevance of Answers

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

**Combination of two factors**

1. Structural: similarity of two nodes in terms of neighbor relationships

2. Distance-based: the PageRank already computed

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Search with Multiple Examples

Combining partial answers



- Multiple Simple Examples

- Each Example describes an Aspect

- Results are Combinations of aspects

- Results have possibly Multiple Structures

Case: Unknown Structures → Find Complex Connections with Simpler Components

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Search Framework

Pruning  and Partial matching



Multi-exemplar Answering

**Input:** Database $G : \langle V, E, \ell \rangle$
**Input:** Samples $\mathcal{S} : \langle s_1, ..., s_m \rangle$
**Output:** Answers $\mathcal{A}$
1: $\mathcal{G} \leftarrow \textsc{Partial}(G, \mathcal{S})$
2: $\mathcal{A} \leftarrow \textsc{Search}(\mathcal{G}, \mathcal{S})$
3: **return** $\mathcal{A}$

Exploit Localized Search

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Fast Candidate Region Search

Reducing the search space

**Identify SEED:**

S1   S3   S2

With cardinality Estimation

Select SINGLE NODE
With neighborhood-mapping

EXPAND around each seed:
Retrieve candidate Regions

DISCARD incomplete regions
With neighborhood-mapping & before graph-search

**SIGIR 2019** tutorial
https://data-exploration.ml

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

https://www.youtube.com/watch?v=A1_dKvX5ZRk

# Graph Query by Example(GQBE)

## Search for example Tuples

**Model:**  Knowledge Graph

**Query:**   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)



Case: Known Entities+Uknown Connections ⟶ Find Complex Connections

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

# GQBE: Maximum Query Graph

Jayaram et al. [2015]

Understand the connections implied by the tuples

$Q = (v_1, v_2)$



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)

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

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# GQBE: Multiple Query Tuples

Understand the connections implied by the tuples

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

**Maximum Query Graph is Very Large**

Subgraphs of Maximum Query graph

GQBE finds answers for multiple query tuples
Compute a re-weighted union graph of the individual query graphs

Preserve the query connectivity

**Full process**

⟨Jerry Yang, Yahoo!⟩

Yahoo! — **founder** → Jerry Yang
Stanford ← **education**
**founded_In** → California
**headquarters** → Sunnyvale
**lived** → Jerry Yang

⟨David Fillo, Yahoo!⟩

⟨Bill Gates, Microsoft⟩

⟨Sergei Brin, Google⟩

Entity Tuple        Maximal Query Graph        Query Lattice        Top-k Answers

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

**SEARCHING** FOR

**Nodes**

**Structures**

BY **LOOKING** AT

| Connectivity | Properties | Queries | (Edge-)Labels |
|---|---|---|---|

**PRODUCES**

Mediator Nodes
[Gionis et al. '15
Ruchansky et al.'15]

Clusters
[Perozzi et al.'14,
Kloumann et al. 14]

Similar Entities
[Metzger et al.'13,
Sobczak et al.'15]

Path Queries
[Bonifati et al.'15]

SPARQL
[Arenas et al.'16]

Entity Tuples
[Jayaram et al.'15]

Similar
Structures
[Mottin et al.'14,
Xie et al.'17,
Lissandrini et al'18]

**Few Approaches accept User Feedback**

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

# Where we are

Relational databases

Textual data

Graphs and networks

Machine learning

Challenges and Remarks

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# How ML fits the Big Picture



**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Interactive exploration of datasets

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

Two ways of learning: passive and active

items

RESPONDS
To USER INPUT

Learn

REQUESTS
USER INPUT

items

Is v ✔ or ✗ ?

Learn

Passive

Active

M. Lissandrini, D. Mottin, 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

M. Lissandrini, D. Mottin, 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

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

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

**SIGIR 2019** tutorial
https://data-exploration.ml

M. Lissandrini, D. Mottin, 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!!!

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

# Explore-by-Example: AIDE

Relevance Feedback

Relevant Samples

Irrelevant Samples

**Data Classification**

User Model

Samples

User Model

Query Formulation

**D-dimensional Space Exploration**

Sampling queries

SQL

Data Extraction Query

**SIGIR 2019** tutorial
https://data-exploration.ml

M. Lissandrini, D. Mottin, 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

**?**

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Classification & Query Formulation

[Dimitriadou et al., 2014, 2016]

| 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

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

# Misclassified Sample Exploitation

[Dimitriadou et al., 2014,2016]



**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Clustering-based Sampling

[Dimitriadou et al., 2014,2016]



Clusters-Sampling Areas

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

**Iterative approach: How many samples does it take to reach the desired result?**

Red wavelength

Green Wavelength

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

# Active learning for online query systems [Vanchinathan et al., 2015]

**Main idea: the system "queries" the user to understand their preferences**



Get item → System → Ask user

preference

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

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Learning unknown preferences

[Vanchinathan et al., 2015]

**Problem**: Find a set S that maximize the unknown 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

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

# A step back …

*Learning from an unknown environment …*

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

# Multi-armed bandits

- Maximize the reward by successively playing gamble machines (the 'arms' of the bandits)

- Invented in **early 1950s** by Robbins for decision making under uncertainty when the environment is unknown

- The reward is unknown ahead of time

Reward $X_1$          Reward $X_2$          Reward $X_3$

**SIGIR 2019** tutorial
https://data-exploration.ml
M. Lissandrini, D. Mottin, T. Palpanas, Y. Velegrakis

# Multi-armed bandits

- Reward = random variable $X_{i,n}$ ; $1 \leq i \leq K, n \geq 1$
- $i$ = index of the gambling machine
- $n$ = number of plays
- $\mu_i$ = expected reward of machine $i$.

A policy, or allocation strategy $A$ is an algorithm that chooses the next machine to play based on the sequence of past plays and obtained rewards.

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

# Exploration vs Exploitation



[https://lilianweng.github.io/lil-log/2018/01/23/the-multi-armed-bandit-problem-and-its-solutions.html](https://lilianweng.github.io/lil-log/2018/01/23/the-multi-armed-bandit-problem-and-its-solutions.html)

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Greedy: A pure exploitation algorithm

Choose the machine with current best expected reward

- **Exploitation vs exploration dilemma**: Should you exploit the information you've learned or explore new options in the hope of greater payoff?

- In the **greedy case**, the balance is completely towards **exploitation**
- Yet, **only exploitation will not lead to a good solution**

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

# Quality measure - Regret

Total expected regret (after T plays):

$$R_T = \mu^* \cdot T - \sum_{i=1}^{K} \mu_j \cdot \mathbb{E}\left[N_{i,T}\right]$$

$\mu^*$: highest expected reward
$\mathbb{E}[N_{i,T}]$: expected number of times machine *i* is played

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# An optimistic view



**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Upper confidence bound (UCB) algorithm

**Optimistic estimate of the mean of arm = 'largest value it could plausibly be'**

1. Pull at each time $t$ the arm with the maximum probability of being the best

$$\frac{1}{n_j} \sum_{s=1}^{n_j} X_{j,s} + \sqrt{\frac{2\log(1/t)}{n_j}}$$

2. Repeat until the budget (number of steps T) is depleted

$n_j$: number of times the arm j has been pulled
**Balance exploration and exploitation:** The uncertainty diminishes as the time passes

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Back to our problem

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

# Modeling the same problem as a Multi-Armed Bandit



Sampling Areas as Arms

Red wavelength

Green Wavelength

**Discrete Search Space Vs Continuous Space?**

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

# Background: Gaussian processes

**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

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

M. Lissandrini, D. Mottin, 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:**
    $[\mu_{t-1}(\cdot), \sigma^2_{t-1}(\cdot)] \leftarrow$ 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

**SIGIR 2019** tutorial
https://data-exploration.ml

M. Lissandrini, D. Mottin, 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

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

Which node to query next?

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

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

# Where could Active learning help?

**Reverse engineering queries and rules**

- Interactive Refinement of example tuples
- Learning the most probable queries from their results

**Graph exploration**

- Summarization of knowledge graphs with preferences
- Seed set expansion
- Recommendation of relevant nodes

**Text processing**

- Fast entity matching
- Advertising based on documents search

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

# Example-based methods

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

# MAB: good resources

**Books and surveys**

- http://slivkins.com/work/MAB-book.pdf

- http://downloads.tor-lattimore.com/book.pdf

- http://sbubeck.com/SurveyBCB12.pdf

**Tutorials**

- Lattimore - AAAI 2018: part 1 - part 2

- Tutorial on bayesian optimization of expensive cost functions

- Blog on bandits: http://banditalgs.com/

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

# Where we are

Relational databases

Textual data

Graphs and networks

Machine learning

**Challenges and Remarks**

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Big data – Easy value?



**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Exploration
*We know where we start*
*we don't know what we'll find*

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

# Traditional Search Methods are not Enough
# We need Specialized Methods for Data Exploration

**From broad views**

**to Detailed view**

**From exploration as**
```
select count(*)
```

**to find what is interesting**

**From Exact Search**

**based on explicit conditions**

**to  Exploratory Search**

**based on Implicit needs**

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Similarities are the key …



**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Example-based methods: All You Need is …



Universe #

Desired Answers $

Examples $\mathcal{E}$

Similarity relation ~

**Implicit (Unknown)**

Query Reverse Engineering

Rule Discovery

Relation Extraction

**Explicit (Known)**

Structural Similarity

Proximity Search

Document Matching

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

# Example-based methods

## Relational

- Reverse engineering queries

- Example-driven schema mapping

- Interactive data repairing

## Textual

- Search documents by example

- Entity extraction by example text

- Web table completion using examples

## Graph

- Community-based Node-retrieval

- Entity Search

- Path and SPARQL queries

- Graph structures as Examples

M. Lissandrini, D. Mottin, 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
- Extract semi-structure data

## Graph

- **Heterogenous Structures**
- Exploit locality
- Entity attributes are expressive
- Large result-sets require ranking

M. Lissandrini, D. Mottin, 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

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

# Acknowledgments

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

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

# Where should we invest time?



**Machine learning**

**Approximate Methods**

**User models**

**Scalability**

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

# Where should we invest time?

**Machine learning**

## Learn from Examples

- … Similarity Measures: are often "fuzzy" and "implicit"

- … New representations of the search space

- Challenge: Scale! Exploration of large search spaces

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

# Where should we invest time?

**Machine learning** → **User models**

## Learn from Examples

- … Similarity Measures to represent User Interests

- … User-centric, dynamic, Exploration-strategies: learn as you go

- Challenge: Distinct User have Different Goals! Explore in different ways

We need more data!

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

# Where should we invest time?



**Scalability**

## Scale Example-based search

- … Huge search space, dynamic data, variety of data models

- … Exploration is Interactive, requires Interactive response time

- Adaptive Data-structures, localized access, flexible schema, incremental index

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

# Where should we invest time?

**Scalability** → **Approximate Methods**

## Scale Example-based search

- … An approximate answer now is better than a precise answer in 1hour

- … Approximate answers can provide insights without being accurate

Exploratory queries retrieve large resultsets: the user needs only a glimpse to  figure out

if they are moving in the right direction!

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# Features of Exploratory Search Systems [White and Roth, 2009]

**Support querying and rapid query refinement:**

- Offer facets and metadata-based <u>result filtering</u>
- Leverage <u>search context</u>

**Example-driven**

- visualizations, summarizations, and <u>explanations</u>
- paired with methods to <u>suggest</u> further example-based explorations.

**Via Interactivity & Personalization!**

**Support learning and understanding**

# Interactive Example Based Exploration System?



**Requires:**

## Fast Query Processing

Avoid the full recomputation of a query

Limit the computation to only a sample

Adaptive query executions

Adaptive data-structures and indexes,

## Automatic Result Analysis

Automatically identify peculiar characteristics,

Data-summarization techniques
Learn user interests automatically

**SIGIR 2019** tutorial
https://data-exploration.ml

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

# ADOPT HETEROGENEITY

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

**operate on heterogeneous datastores**

**dataset search**

*Data Lakes??*

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

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

Tools that work on **commodity hardware, mobile devices**

Data-exploration for **everyday use-cases**

Users want back **the control on their data**

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

# NATARUAL LANGUAGE INTERFACE

*flexible, vague, imprecise input*

**Exploration through conversation**

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

# Example is always more efficacious than precept

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

**"New Trends on Exploratory Methods for Data Analytics."** *PVLDB, 2017.*

**"Data Exploration Using Example-Based Methods."** *M&C, 2018.*

**"Exploring the Data Wilderness through Examples."** *SIGMOD, 2019.*

**Slides:** https://data-exploration.ml/

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

# 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.

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

# 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.

He, J., Veltri, E., Santoro, D., Li, G., Mecca, G., Papotti, P. and Tang, N. Interactive and deterministic data cleaning. SIGMOD, 2016.

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

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

# References

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

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

M. Lissandrini, D. Mottin, Y. Velegrakis, T. Palpanas. Multi-Example Search in Rich Information Graphs ICDE 2018

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 Sanchez, and E. Müller. Focused clustering and outlier detection in large attributed graphs. KDD, 2014.

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

# 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.

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

# 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.

**SIGIR 2019** tutorial
https://data-exploration.ml
M. Lissandrini, D. Mottin, T. Palpanas, Y. Velegrakis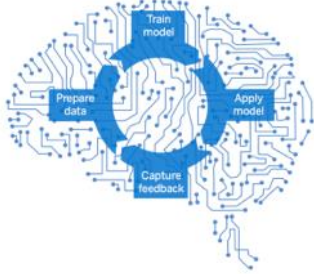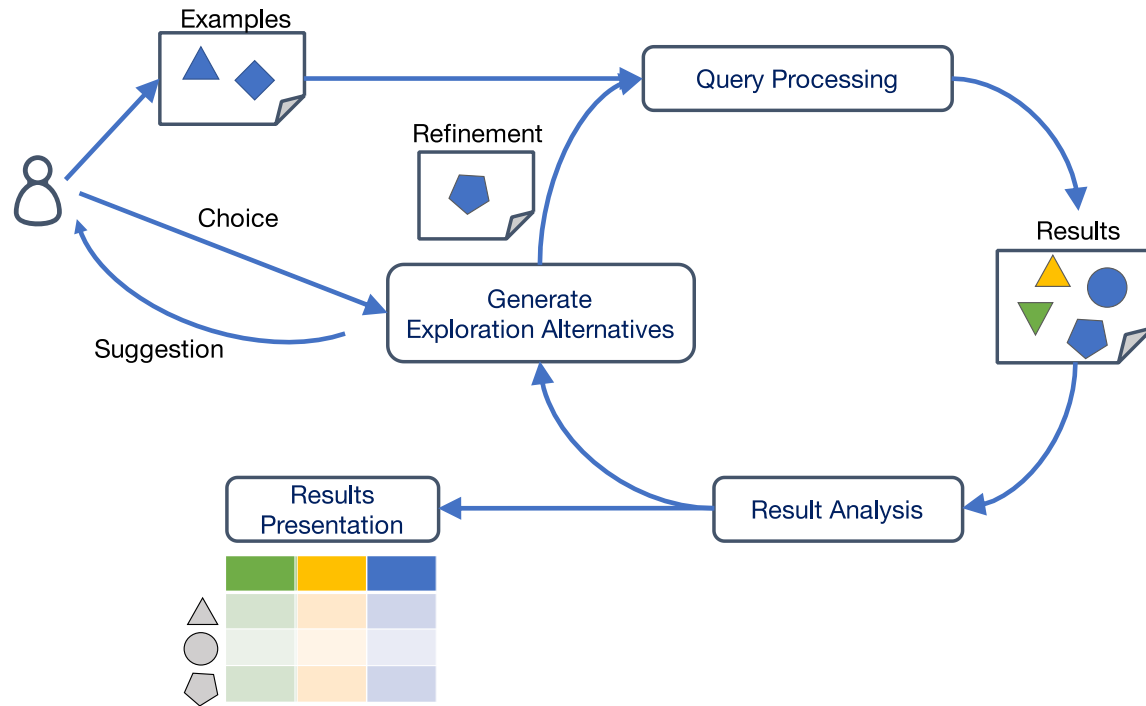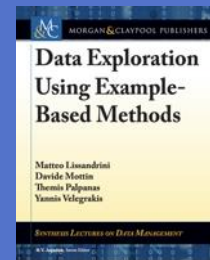