# Two decades of Pattern Mining eBISS2016, Tours

*arnaud.soulet@univ-tours.fr

UNIVERSITÉ FRANÇOIS - RABELAIS TOURS

LI

**Outline**

**What is Pattern Mining?**

20 years ago...
*SIGMOD Conference 1993

**What are the main principles?**

Typology of publications
❶ Language
❷ Constraint
❸ Condensed Representation (CR)
*Levelwise Search and Borders of Theories in Knowledge Discovery
[Manila and Toivonen, 1997]

**What are the recent trends?**

New trends
*Recent keywords of Pattern Mining

# What is Pattern Mining?

# Mining Association Rules between Sets of Items in Large Databases

20 years ago...

*SIGMOD Conference 1993

**New problem**

Mining Association Rules between Sets of Items in Large Databases

We introduced the problem of mining association rules between sets of items in a large database of customer transactions. Each transaction consists of items purchased by a customer in a visit. We are interested in finding those rules that have:

- Minimum transactional support $s$ — the union of items in the consequent and antecedent of the rule is present in a minimum of $s\%$ of transactions in the database.

- Minimum confidence $c$ — at least $c\%$ of transactions in the database that satisfy the antecedent of the rule also satisfy the consequent of the rule.

monthly purchases by members of a book club or a music club. Several organizations have collected massive amounts of such data. These data sets are usually stored

and "mustard" in the consequent. This query can be phrased alternatively as a request for the additional items that have to be sold together with sausage in order to make it highly likely that mustard will also

An example of such an association rule is the statement that 90% of transactions that purchase bread and butter also purchase milk. The antecedent of this rule consists of bread and butter and the consequent consists of milk alone. The number 90% is the confidence factor of the rule.

**\*Discovering all relevant association rules**

**New solution**

We introduced the problem of mining association rules between sets of items in a large database of customer transactions. Each transaction consists of items purchased by a customer in a visit. We are interested in finding those rules that have:

- Minimum transactional support $s$ — the union of items in the consequent and antecedent of the rule is present in a minimum of $s\%$ of transactions in the database.

- Minimum confidence $c$ — at least $c\%$ of transactions in the database that satisfy the antecedent of the rule also satisfy the consequent of the rule.

monthly purchases by members of a book club or a music club. Several organizations have collected massive amounts and "mustard" in the consequent. This query can be phrased alternatively as a request for the additional items that have to be sold together with sausage is

We solve this problem by decomposing it into two subproblems:

1. Finding all itemsets, called *large* itemsets, that are present in at least $s\%$ of transactions.

2. Generating from each large itemset, rules that use items from the large itemset.

**\*Enumerating all frequent itemsets**

**Leonardo DiCaprio Plays the Same Character Over and Over**

|  | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| Titanic | ✓ | | ✓ | |
| Catch Me If You Can | ✓ | | | ✓ |
| Inception | ✓ | ✓ | ✓ | |
| Django | | ✓ | ✓ | |
| The Great Gatsby | ✓ | ✓ | ✓ | ✓ |

**Troubled romantic + Rich**

**(supp = 0.4)**

*Itemset

**Leonardo DiCaprio Plays the Same Character Over and Over**

|  | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| *Titanic* | ✓ |  | ✓ |  |
| *Catch Me If You Can* | ✓ |  |  | ✓ |
| *Inception* | ✓ | ✓ | ✓ |  |
| *Django* |  | ✓ | ✓ |  |
| *The Great Gatsby* | ✓ | ✓ | ✓ | ✓ |

Troubled romantic ➔ Rich
(supp = 0.4 / conf = 0.5)

Troubled romantic ➔ Dies
(supp = 0.6 / conf = 0.75)

**Leonardo DiCaprio Plays the Same Character Over and Over**

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| Titanic | A | | C | |
| Catch Me If You Can | A | | | D |
| Inception | A | B | C | |
| Django | | B | C | |
| The Great Gatsby | A | B | C | D |

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

# Frequent patterns

Ø

A    B    C    D

AC    AB    BC    BD    AD    CD

ABC    ABD    ACD    BCD

ABCD

*Minimal frequency threshold = 1

| Troubled Romantic | Rich | Dies | Hidding Secret |
|:---:|:---:|:---:|:---:|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

# Frequent patterns

```
                    Ø

        A       B       C       D

AC    AB    BC    BD    AD    CD

    ABC   ABD   ACD   BCD

        ABCD
```

**\*Minimal frequency threshold = 2**

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

# Frequent patterns

Ø

A    B    C    D

AC   AB   BC   BD   AD   CD

ABC   ABD   ACD   BCD

ABCD

*Minimal frequency threshold = 3

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A |  | C |  |
| A |  |  | D |
| A | B | C |  |
|  | B | C |  |
| A | B | C | D |

# Frequent patterns

Ø

| A | B | C | D |
|---|---|---|---|

AC   AB   BC   BD   AD   CD

ABC   ABD   ACD   BCD

ABCD

*Minimal frequency threshold = 4

Exact solution

Exhaustive search

Speed of answer

Pattern Mining

*The footprint of databases

**Exact solution**

**Approximate solution**

**Exhaustive search**

**Heuristic search**

**Speed of answer**

**Quality of solution**

**Pattern Mining vs Artificial Intelligence**

*The footprint of databases

**Rakesh Agrawal**

" I'm a **database** person, so my view of data mining has been that it is essentially a richer form of querying."

*The footprint of databases

# Impact of this seminal paper?

**Classical survey**

**55 references**

**45 references**

***Dozens of references**

**Classical survey**

*Dozens of references

**Bibliometric survey**

*Thousands of references

# Materials

**KDD**
**1994**

**PAKDD**
**PKDD**
**1997**

**ICDM**
**SDM**
**2001**

**1993**

*Data mining conferences ranked A

# Materials



[Agrawal and Srikant, 1994]

*Not consider VLDB, CIKM, ICDE,...

**Materials**

KDD 199...

ICDM
SDM
2001

[Manila and Toivonen, 1997]

*Not consider DMKD, TKDE, KAIS,...

KDD :      1,905 since 1995

PKDD :     1,295 since 1997

PAKDD :    1,277 since 1998

ICDM :     1,598 since 2001

Materials  SDM :   813 since 2002

*6,888 publications from DBLP (1995-2012)

**Faster, Higher, Stronger**

Pattern or not?  Fuzzy limit

*Discovery and use

❶ **Keyword filtering** for selecting good candidate papers

**Pattern or not?** ❷ **Manual filtering** for removing False Positive

*Semi-automated topic assignment

Language

Constraint

LOCAL

GLOBAL

Pattern or not?  Condensed Representation

*Dimensions of Pattern Mining

pattern, item, sequence, rule, tree, graph, string, stream, subgroup...

support (no Vector Machine), frequent, monotone...

LOCAL

GLOBAL

Pattern or not? free, generator, closed, condensed, concise

*Keywords of Pattern Mining

Keyword **filtering**             1,732

Manual **filtering**              1,087

**Pattern or not?**          (148 **consulted abstracts**)

*****5% of False Negative, around 258 papers**

1/5th of authors and 1/6th of KDD publications

PM; 1087

Other papers; 5801

PM; 1789

Other authors; 7579

*PM is a true subfield of KDD

The slowdown of Pattern Mining

Other

PM

1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012

*Turning point: 2006

The slowdown of Pattern Mining

*Golden age: 1998-2005 (1 paper out of 5)

# What are the main principles of Pattern Mining?

**Typology of publications**

❶ **Language**

❷ **Constraint**

❸ **Condensed Representation (CR)**

*Levelwise Search and Borders of Theories in Knowledge Discovery

[Manila and Toivonen, 1997]

## Levelwise Search and Borders of Theories in Knowledge Discovery

HEIKKI MANNILA                                              heikki.mannila@cs.helsinki.fi
HANNU TOIVONEN                                            hannu.toivonen@cs.helsinki.fi
Department of Computer Science, P.O. Box 26, FIN-00014 University of Helsinki, Finland

**Abstract.** One of the basic problems in knowledge discovery in databases (KDD) is the following: given a data set $r$, a class $\mathcal{L}$ of sentences for defining subgroups of $r$, and a selection predicate, find all sentences of $\mathcal{L}$ deemed interesting by the selection predicate. We analyze the simple levelwise algorithm for finding all such descriptions. We give bounds for the number of database accesses that the algorithm makes. For this, we introduce the concept of the border of a theory, a notion that turns out to be surprisingly powerful in analyzing the algorithm. We also consider the verification problem of a KDD process: given $r$ and a set of sentences $S \subseteq \mathcal{L}$, determine whether $S$ is exactly the set of interesting statements about $r$. We show strong connections between the verification problem and the hypergraph transversal problems. The verification problem arises in a natural way when using sampling to speed up the pattern discovery step in KDD.

**Keywords:** theory of knowledge discovery, association rules, episodes, integrity constraints, hypergraph transversals
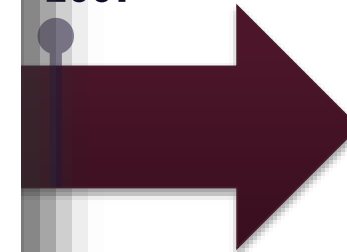
### 1. Introduction

Knowledge discovery in databases (KDD), also called data mining, has recently received wide attention from practitioners and researchers. There are several attractive application areas for KDD, and it seems that techniques from machine learning, statistics, and databases can be profitably combined to obtain useful methods and systems for KDD. See, e.g., Fayyad et al. (1996), and Piatetsky-Shapiro and Frawley (1991) for general descriptions of the area.

The KDD area is and should be largely guided by (successful) applications. Still, theoretical work in the area is needed. In this paper we take some steps towards theoretical KDD. We consider a KDD process in which the analyzer first produces lots of potentially interesting rules, subgroup descriptions, patterns, etc., and then interactively selects the truly interesting ones from these. In this paper we analyze the first stage of this process: how to find all the potentially interesting rules in the database.

The intuitive idea behind this work is as follows. A lot of work in data mining can be formulated in terms of finding all rules of the form

if $\varphi$ then $\vartheta$,

The World of Pattern Mining

**1**

**Language?**

***Do you remember?**

**1**

Language?

The World of Pattern Mining

## Items, itemsets and datasets

- **Items = a set of distinct literals**
  - Example:
  - Often denoted by a letter of alphabet
- **Itemset = any set of items**
  - Example:
  - The whole set of itemsets is called the **language**
- **Dataset = multi-set of itemsets**
  - Example:

27/06/2016　　　　Data mining - Local pattern discovery　　　　19

\*Do you remember?

# The World of Pattern Mining

## Items, itemsets and ...

## Formal Notations of Sequence Data

- An *item* is denoted by lower-case letters

$$a, b, c, d, \ldots$$

- An *itemset* is denoted by upper-case letters

$$I = (abc), \quad I' = (acd), \quad I_1 = (bcef), \ldots$$

- A *sequence* is denoted by $s$ (with prime and/or indice)

$$s = \langle I_1 I_2 I_3 I_4 \rangle, \quad s' = \langle (ab)(c)(bcd)(ac) \rangle, \quad s_1 = \langle (ab)c(bc) \rangle \ldots$$

(We may ignore the parentheses for 1-item itemsets)

- A *sequence database* is a (large) set of sequences

**1**

**Language?**

**ber?**

8 itemsets

∅,A, B, C, AB, AC, BC,ABC

80 sequential patterns

∅,<A>,< AA>, <AAA>, <AAB>, <AAC>, ...

**Language sophistication with 3 items and 3 as maximal length**

238 subgraphs patterns

∅,A,AA,A-A,AAA,A-AA,A-A-A,...

*Pattern explosion

Pattern explosion

Pattern matching*

Computational challenges of language sophistication

Subgraph isomorphism checking

*Does a database entry contain a pattern?

Pattern explosion

Pattern matching

Computational challenges of language sophistication

Subgraph isomorphism checking*

*Does a graph contain a subgraph isomorphic to another graph?

itemset, set

rule, association

sequence, episode, string, stream, protein, periodic, temporal

graph, molecular, structure, network

tree, xml

**Keywords about language** spatial, spatio-temporal

relational

*Semi-automated topic assignment

**Itemsets or rules for 63% of publications**

rule
32%

itemset
31%

sequence
17%

graph
10%

tree
4%

spatial
3%

generic
2%

relational
1%

*Only 18 papers with generic language

**Sophistication and marginalization of languages**

*Next language: spatio-temporal patterns?

# Progress of sequences and graphs

Rule
Itemset
Sequence
Graph

80%
70%
60%
50%
40%
30%
20%
10%
0%

1993-1996   1997-2000   2001-2004   2005-2008   2009-2012

**Next challenge: complex data**

*Uncertain, dynamic, massive, heterogeneous data

**Pattern explosion of frequent patterns (even with itemsets)**



*How to reduce the number of patterns?

**Constraint**

Focusing on the most useful patterns for the data expert

**Condensed Representation**

Removing all redundant patterns

Two strategies against pattern explosion

[Vreeken et al., 2010]

*Useful Patterns (UP'10) ACM SIGKDD Workshop

## Support and confidence

- **Support of X** = proportion of transactions in dataset D containing X

$$supp(X, D) = |\{t \in D / X \subseteq t\}| / |D|$$

- **Support of X→Y**

$$supp(X \rightarrow Y, D) = supp(X \cup Y, D)$$

- **Confidence of X→Y**

$$conf(X \rightarrow Y, D) = supp(X \cup Y, D)/supp(X, D)$$

27/06/2016

Data mining - Local pattern discovery

25

**Constraint?**

**\*Do you remember, again?**

**❷**

**Constraint?**

*Do you remember, again?*

## Emerging patterns: definition

- Growth rate of X in $D_i$:
  $gr_i(X,D)=supp(X,D_i)/supp(X,D-D_i)$
- X is an emerging pattern iff $gr_i(X,D) \geq \rho$   $(\rho>1)$

Class 1 / Class 2 — $gr_1(X,D)>1$

Class 1 / Class 2 — $gr_1(X,D)<1$

Class 1 / Class 2 — $gr_1(X,D)=1$

Efficient Mining of Emerging Patterns: Discovering Trends and Differences. Dong et Li, KDD 1999.

Data mining - Local pattern discovery

27/06/2016

27/06/2016

78

25

# ❷

**Constraint?**

## Downward and upward closure

- S
- ...
- S
- C

specialization →

Itemsets satisfying an anti-monotone constraint

Itemsets satisfying a monotone constraint

$freq(X, D) \geq t$
$min(X.val) \geq t$
$max(X.val) \leq t$
$sum(X.val) \leq t$
$X \subseteq \{A, B, C\}$
...

$freq(X, D) \leq t$
$min(X.val) \leq t$
$max(X.val) \geq t$
$sum(X.val) \geq t$
$X \supseteq \{A, B, C\}$
...

27/06/201

26/10/2011

Data mining – Local pattern discovery

84

**Do you remember, again?**

|  | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|

Troubled romantic ➔ Rich (conf = 0.5)

Interest of constraints

Troubled romantic ➔ Dies (conf = 0.75)

|  | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| Titanic | ✓ | | ✓ | |
| Catch Me If You Can | ✓ | | | ✓ |
| Inception | ✓ | ✓ | ✓ (spinning top) | |
| Django | | ✓ | ✓ | |
| The Great Gatsby | ✓ | ✓ | ✓ | ✓ |

**Interest of constraints**

Troubled romantic ➔ Rich

(conf = 0.5 / lift = 0,8)

Troubled romantic ➔ Dies

(conf = 0.75 / lift = 0,6)

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

$$freq(X) \geq minfreq$$

Ø

A    B    C    D

AC   AB   BC   BD   AD   CD

ABC   ABD   ACD   BCD

ABCD

**Challenge of constraints**

**\*How to prune the search space for frequency?**

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

$$freq(X) \geq minfreq$$

Ø

A  B  C  D

AC  AB  BC  BD  AD  CD

ABC  ABD  ACD  BCD

ABCD

**Challenge of constraints**

*How to prune the search space for frequency?
Easy due to the downard closure

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

$$freq(X) \times |X| \geq minarea$$

Ø

A    B    C    D

AC  AB  BC  BD  AD  CD

ABC  ABD  ACD  BCD

ABCD

**Challenge of constraints**

**\*How to prune the search space for area?**

|  | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
|  | A |  | C |  |
|  | A |  |  | D |
|  | A | B | C |  |
|  |  | B | C |  |
|  | A | B | C | D |

$$freq(X) \times |X| \geq minarea$$

∅

A  B  C  D

AC  AB  BC  BD  AD  CD

ABC  ABD  ACD  BCD

ABCD

**Challenge of constraints**

***How to prune the search space for area?***

**Relax the area constraint by** $freq(X) \times 4 \geq minarea$

regularity, frequent, support

contrast, emerging, discriminative

exception, abnormal, surprising, anomaly, unexpected

utility

significant, chi-square, correlated

interesting, relevant

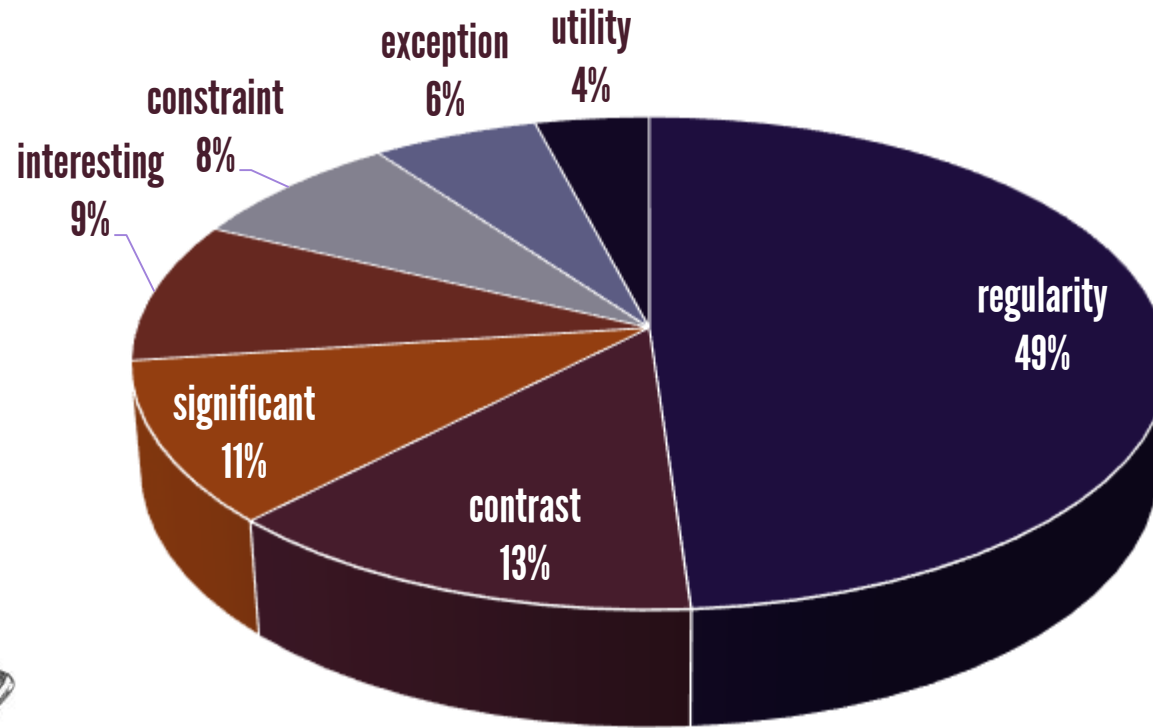**Keywords about constraints** generic, monotone, anti-monone, constrained
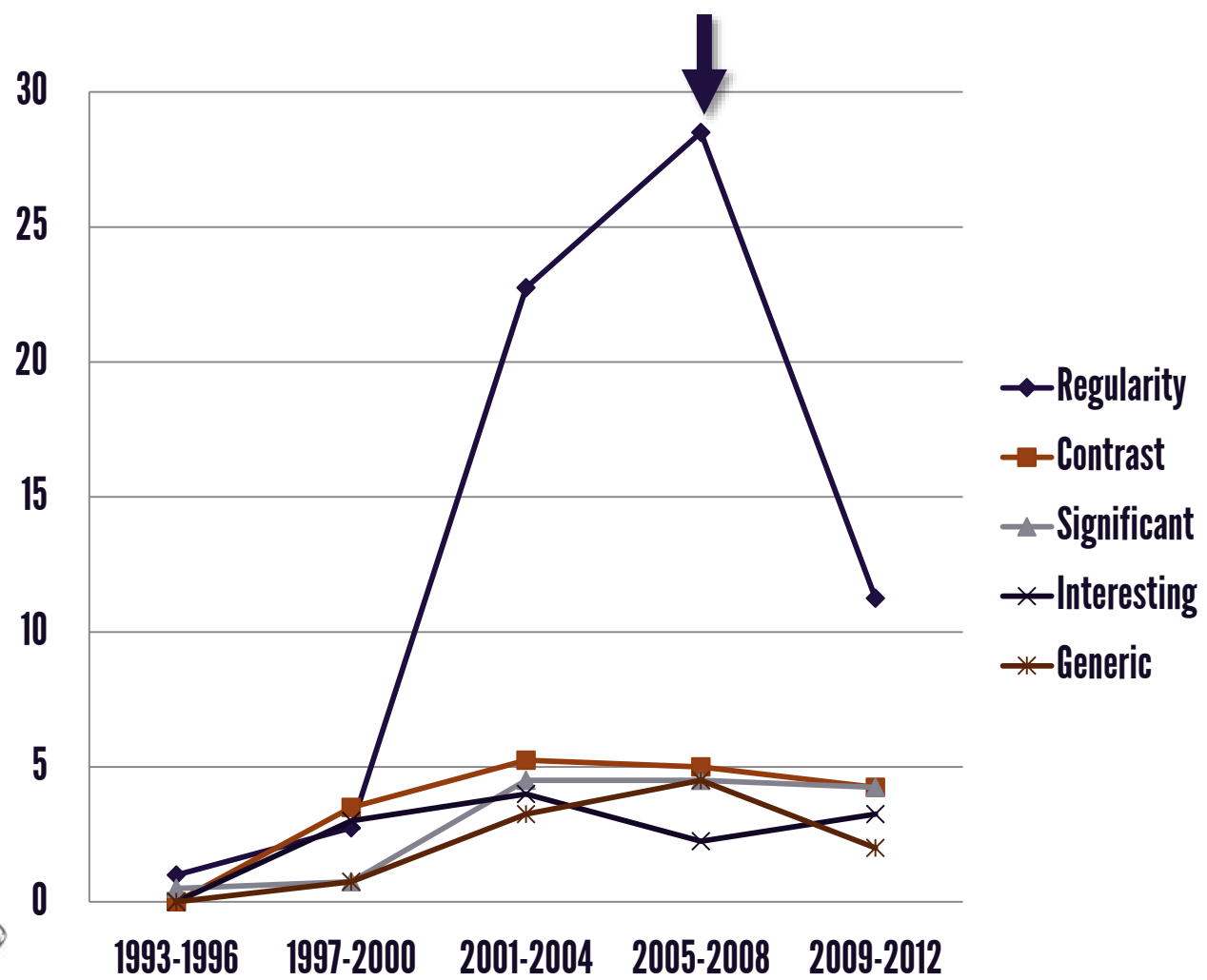
**Frequent patterns in 49% of publications**

regularity
49%

contrast
13%

significant
11%

interesting
9%

constraint
8%

exception
6%

utility
4%
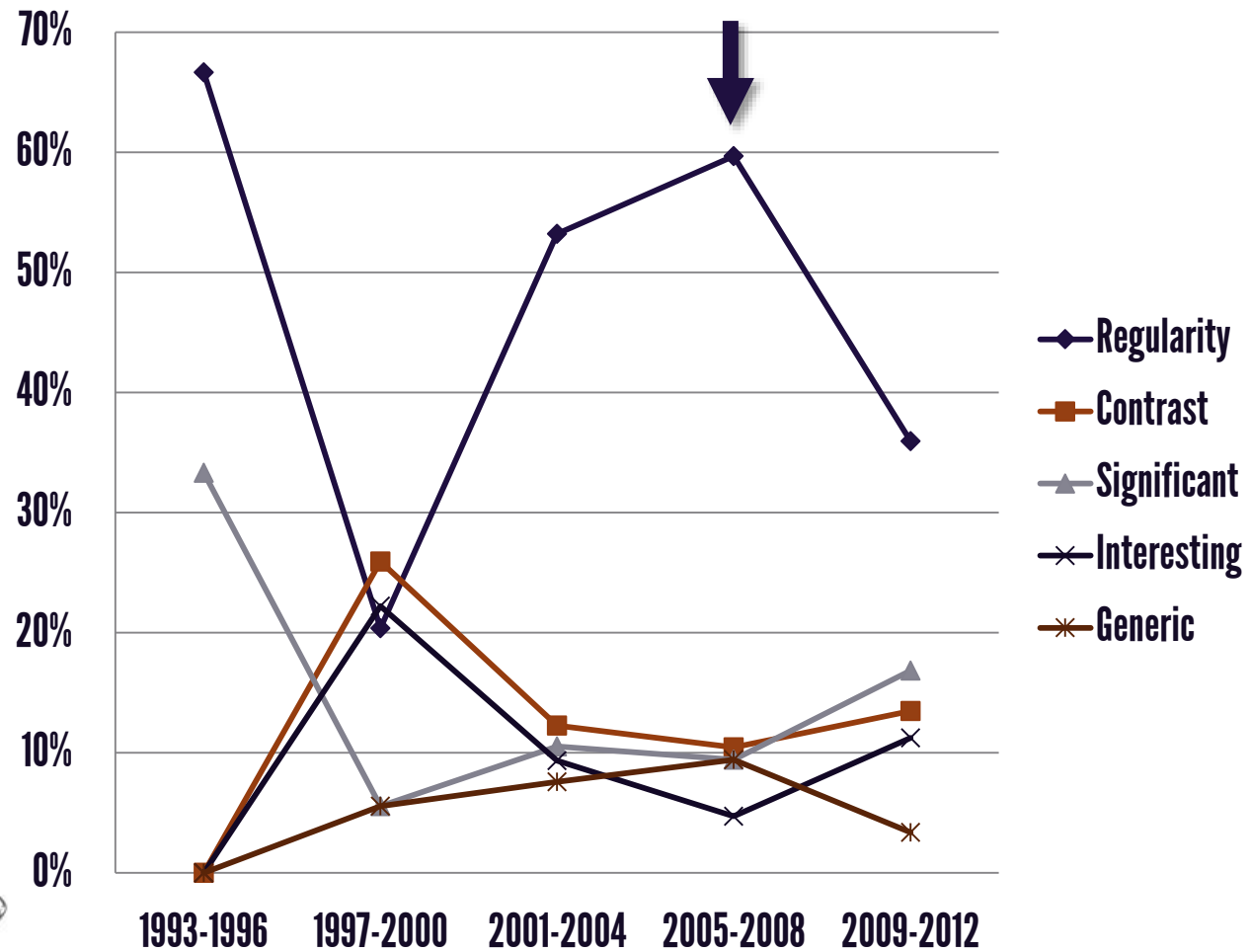
*Speed preferred to quality

**Stable except for frequency**

# Progress of contrastive and significant patterns



*Only 42 papers with generic constraints

# Interestingness

Piatetsky-Shapiro
Agrawal
2003

Han et al.
2007

"we need work to bring in some notion of 'here is my idea of what is interesting,' and pruning the generated rules

based on that input."

Rakesh Agrawal Speaks Out
on Where the Data Mining Field Is Going, Where It Came From, How to Choose Problems and Open Up New Fields, Our Responsibilities to Society as Technologists, What Industry Owes Academia, and More

by Marianne Winslett

# Interestingness

Piatetsky-Shapiro
Agrawal
2003

Han et al.
2007

**Frequent pattern mining: current status and future directions**

Jiawei Han · Hong Cheng · Dong Xin · Xifeng Yan

**Abstract** Frequent pattern mining has been a focused theme in data mining research for over a decade. Abundant literature has been dedicated to this research and tremendous progress has been made, ranging from efficient and scalable algorithms for frequent itemset mining in transaction databases to numerous research frontiers, such as sequential pattern mining, structured pattern mining, correlation mining, associative classification, and frequent pattern-based clustering, as well as their broad applications. In this article, we provide a brief overview of the current status of frequent pattern mining and discuss a few promising research directions. We believe that frequent pattern mining research has substantially broadened the scope of data analysis and will have deep impact on data mining methodologies and applications in the long run. However, there are still some challenging research issues that need to be solved before frequent pattern mining can claim a cornerstone approach in data mining applications.

**Keywords** Frequent pattern mining · Association rules · Data mining research · Applications.

"it is still not clear what kind of patterns will give us satisfactory pattern sets in both compactness and representative quality"

# ❸

## Condensed representations



*Do you remember, again and again?

**③**

**Condensed representations**

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

**Maximal patterns**

Ø

A   B   C   D

AC  AB  BC  BD  AD  CD

ABC  ABD  ACD  BCD

ABCD

**2 patterns with minfreq = 2**

***larger frequent patterns w.r.t inclusion**

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

**Closed patterns**

**7 patterns with minfreq = 2**

**\*maximal patterns of equivalence classes**

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

**Free patterns**

Ø

A    B    C    D

AC    AB    BC    BD    AD    CD

ABC    ABD    ACD    BCD

ABCD

**7 patterns with minfreq = 2**

**\*minimal patterns of equivalence classes**

# More condensed representations

**Non-derivable itemsets 2002**

**K-free itemsets 2003**

Mining All Non-Derivable Frequent Itemsets

[Calders and Goethals, 2002]

Minimal k-Free Representations of Frequent Sets

[Calders and Goethals, 2003]

**Keywords about condensed representations**

closed

border, maximal, minimal

free, generator, non-derivable

*Semi-automated topic assignment

**Closed patterns in 60% of CR**

concision 7%

free 13%

border 20%

closed 60%

*Why this success?

**2005-2008:**
**Activity peak of CRs**

*Pattern-based models instead of CRs

**Progress of free itemsets**

# Pattern-based classifiers

CBA
1998

*Two-phases: pattern extraction & model construction

# Pattern-based models

**Krimp**

2006

**The Chosen Few**

**MINI**

2007

[Bringmann and Zimmermann, 2007]

[Vreeken et al., 2006]

[Gallo et al., 2007]

*Two-phases: pattern extraction & model construction

# Pattern-based models

MUSK
2009

Local pattern
sampling
2011

[Hasam and Zaki, 2009]

[Boley et al., 2011]

*Sampling

**Exact solution**      **Approximate solution**

**Exhaustive search**      **Heuristic search**

**Speed of answer**      **Quality of solution**

**Pattern Mining vs Artificial Intelligence**

*The footprint of databases

**Pattern-based models**

MUSK 2009

Local pattern sampling 2011

**Stop completeness!**

Useful Patterns ACM SIGKDD Workshop 2010

ECMLPKDD most-influential paper award 2012

"Please, please stop making new algorithms for mining all patterns"

Toon Calders

1993

now*

*Aprroximate solution and Heuristic search!

# What are the recent trends of Pattern Mining?

# New trends



*Recent keywords of Pattern Mining

# New trends



*Complex data

# New trends



*From speed to quality

# Pattern mining as an optimization problem

**Top-k frequent patterns**
**2000**

**Skypatterns**
**2011**

**Optimal patterns**
**Dominance programming**
**2015**



Mining N-most Interesting Itemsets

[Fu et al., 2000]

Mining Dominant Patterns in the Sky

[Soulet et al., 2011]

Modeling and Mining Optimal Patterns using Dynamic CSP

[Urgate et al., 2015]

*Focusing on the best patterns

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

Ø

A        B        C        D

AC     AB     BC     BD     AD     CD

ABC    ABD    ACD    BCD

ABCD

**Top-k pattern mining**

**\*Finding the k patterns maximizing an interestingness measure**

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

Top-k pattern mining

Ø

A    B    C    D

AC    AB    BC    BD    AD    CD

ABC    ABD    ACD    BCD

ABCD

*Finding the 3 most frequent patterns: Ø (5), A (4), C (4)

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

**Top-k pattern mining**

Ø

A    B    C    D

AC   AB   BC   BD   AD   CD

ABC   ABD   ACD   BCD

ABCD

*Finding the 3 most frequent patterns: Ø (5), A (4), C (4)

**Easy due to anti-monotone property of frequency

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

**Top-k pattern mining**

Ø

A　　B　　C　　D

AC　AB　BC　BD　AD　CD

ABC　ABD　ACD　BCD

ABCD

*Finding the 3 patterns maximizing area: AC (6), BC (6), ABC (6)

**Branch&Bound method

**Top-k pattern mining**

Compact

Threshold free

Best patterns

Not fast*

No diversity

*Exact resolution is costly / sometimes heuristic search (beam seacrh)

**Top-k pattern mining**

Compact

Threshold free

Best patterns

Not fast

No diversity*

*Diversity issue: top-k patterns often very similar

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

**Skyline pattern mining**

| Pattern | Freq. | Area | |
|---|---|---|---|
| Ø | **5** | 0 | Top frequent |
| A | 4 | 4 | |
| C | 4 | 4 | |
| AC | 3 | **6** | Top area |
| BC | 3 | **6** | |
| AD | 2 | 4 | |
| ABC | 2 | 6 | |
| ABCD | 1 | 4 | |

**\*How to find a trade-off between several criteria?**

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

**Skyline pattern mining**

| Pattern | Freq. | Area |
|---|---|---|
| Ø | 5 | 0 |
| A | 4 | 4 |
| C | 4 | 4 |
| AC | 3 | 6 |
| BC | 3 | 6 |
| AD | 2 | 4 |
| ABC | 2 | 6 |
| ABCD | 1 | 4 |

Skyline pattern mining

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

| Pattern | Freq. | Area |
|---|---|---|
| Ø | 5 | 0 |
| A | 4 | 4 |
| C | 4 | 4 |
| AC | 3 | 6 |
| BC | 3 | 6 |
| AD | 2 | 4 |
| ABC | 2 | 6 |
| ABCD | 1 | 4 |

Area

ABC    AC, BC

ABCD    AD    A, C

Ø    Freq.

*Dominated space

# Skyline pattern mining

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

| Pattern | Freq. | Area |
|---|---|---|
| Ø | 5 | 0 |
| A | 4 | 4 |
| C | 4 | 4 |
| AC | 3 | 6 |
| BC | 3 | 6 |
| AD | 2 | 4 |
| ABC | 2 | 6 |
| ABCD | 1 | 4 |

Area — Top area AC, BC — A, C — Ø Top frequent — Freq.

**\*Skypatterns = non-dominated patterns**

Skyline pattern mining

| Pattern | Freq. | Area |
|---------|-------|------|
| ∅ | 5 | 0 |
| A | 4 | 4 |
| C | 4 | 4 |
| AC | 3 | 6 |
| BC | 3 | 6 |
| AD | 2 | 4 |
| ABC | 2 | 6 |
| ABCD | 1 | 4 |

*Skypatterns are closed patterns

Maximal patterns

Closed patterns

Dominance programming for optimal patterns

Top-k patterns

Skypatterns

*A pattern is optimal if it is not dominated by another.

# Maximal patterns*

# Closed patterns

# Dominance programming for optimal patterns

# Top-k patterns

# Skypatterns

*Dominance relation = inclusion

Dominance programming for optimal patterns

Maximal patterns

Closed patterns*

Top-k patterns

Skypatterns

*Dominance relation = inclusion at same frequency

Maximal patterns

Closed patterns

Dominance programming for optimal patterns

Top-k patterns*

Skypatterns

*Dominance relation = order induced by the interestingness measure

**Dominance programming for optimal patterns**

**Maximal patterns**

**Closed patterns**

**Top-k patterns**

**Skypatterns***

***Dominance relation = Pareto domination**

**Dominance programming for optimal patterns**

Closed patterns

Maximal patterns

Skypatterns

Top-k patterns

# New trends



*From exhaustive collection to models

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

ø

A    B    C    D

AC    AB    BC    BD    AD    CD

ABC    ABD    ACD    BCD

ABCD

**Pattern sampling**

**Pattern A (freq = 4) has twice more chance to be drawn than pattern D (freq = 2)**

***Picking k patterns randomly with a probability proportional to an interestingness measure**

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

**Pattern sampling**

**Stochastic methods** [Hasan and Zaki, 2009]

**Random walk on lattice**

**Two-step direct method** [Boley et al., 2011]

**Pick a transaction + pick an itemset of this transaction**

***Two main families of methods**

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

**Pattern sampling**

**Stochastic methods** [Hasan and Zaki, 2009]

**Random walk on lattice**

**Two-step direct method** [Boley et al., 2011]*

**Pick a transaction + pick an itemset of this transaction**

***More uniform**

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

Ø, A, C, AC

Ø, A, D, AD

Ø, A, B, C, AB, AC, BC, ABC

Ø, B, C, BC

Ø, A, B, C, D, AB, AC, AD, BC, BD,
CD, ABC, ABD, ACD, BCD, ABCD

# Direct pattern sampling

[Boley et al., 2011]

*Consider all itemsets contain in each transaction

|  | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

ø, **A**, C, AC

ø, **A**, **D**, AD

ø, **A**, B, C, AB, AC, BC, ABC

ø, B, C, BC

ø, **A**, B, C, **D**, AB, AC, AD, BC, BD, CD, ABC, ABD, ACD, BCD, ABCD

**Direct pattern sampling**

[Boley et al., 2011]

Pattern **A** (freq = 4) appears twice more than pattern **D** (freq = 2)

**\*Consider all itemsets contain in each transaction**

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

**Direct pattern sampling**

[Boley et al., 2011]

4    ø, **A**, C, AC

4    ø, **A**, **D**, AD

8    ø, **A**, B, C, AB, AC, BC, ABC

4    ø, B, C, BC

16   ø, **A**, B, C, **D**, AB, AC, AD, BC, BD, CD, ABC, ABD, ACD, BCD, ABCD

Pattern **A** (freq = 4) has twice more chance to be drawn than pattern **D** (freq = 2)

*Count the number of itemsets

| | Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|---|
| | A | | C | |
| | A | | | D |
| | A | B | C | |
| | | B | C | |
| | A | B | C | D |

**Direct pattern sampling**

[Boley et al., 2011]

1/9  ∅, A, C, AC

1/9  ∅, A, D, AD

2/9  ∅, A, B, C, AB, AC, BC, ABC

1/9  ∅, B, C, BC

4/9  ∅, A, B, C, D, AB, AC, AD, BC, BD, CD, ABC, ABD, ACD, BCD, ABCD

Pattern A (freq = 4) has twice more chance to be drawn than pattern D (freq = 2)

*Normalize

| Troubled Romantic | Rich | Dies | Hidding Secret |
|---|---|---|---|
| A | | C | |
| A | | | D |
| A | B | C | |
| | B | C | |
| A | B | C | D |

**Direct pattern sampling**

[Boley et al., 2011]

1/9    Ø, **A**, C, AC

1/9    Ø, **A**, **D**, AD

2/9    Ø, **A**, B, C, AB, AC, BC, ABC

1/9    Ø, B, C, BC

4/9    Ø, **A**, B, C, **D**, AB, AC, AD, BC, BD, CD, ABC, ABD, ACD, BCD, ABCD

**Pattern A (freq = 4) has twice more chance to be drawn than pattern D (freq = 2)**

*Pick a transaction proportionally to the distribution

**Pick uniformly an itemset within this transaction

**Pattern sampling**

Compact

Threshold free

Diversity

Very fast

Patterns far from optimality

**Ease of use**

**Constraint based
pattern mining**

**\*No algorithm specification**

Ease of use

Constraint based
pattern mining

Pattern sampling*

Optimal pattern mining*

*No user-specified threshold

Ease of use

Constraint based
pattern mining

Pattern sampling

Optimal pattern mining

Interactive
pattern mining*

*No user-specified measure

**Interactive data exploration using pattern mining**

[van Leeuwen 2014]

Mine

Interact

Learn

**Interactive data exploration using pattern mining**

[van Leeuwen 2014]

Mine

Candidate patterns

Feedback integration

Interact

Learn

User's feedback

**Interactive data exploration using pattern mining**

[van Leeuwen 2014]

**Mine**

Feedback integration

Candidate patterns*

**Learn**

**Interact**

User's feedback

*Active learning vs useful pattern mining

Interactive data exploration using pattern mining
[van Leeuwen 2014]

Mine

Feedback integration

Candidate patterns

Learn

Interact

User's feedback*

*Explicit feedback vs implicit feedback

Interactive data exploration using pattern mining

[van Leeuwen 2014]

Mine

Feedback integration*

Candidate patterns

Learn

Interact

User's feedback

*How to upate the target of the mining method?

Discovering Interesting Patterns Through User's Interactive Feedback [Xin et al., 2006]*

Interactive Pattern Mining on Hidden Data: A Sampling-based Solution [Bhuiyan et al., 2012] **

Active Preference Learning for Ranking Patterns [Dzyuba et al., 2013] **

Mining step?

*Offline mining of all frequent patterns
**Online mining by integrating preferences

Discovering Interesting Patterns Through User's Interactive Feedback [Xin et al., 2006]

Interactive Pattern Mining on Hidden Data: A Sampling-based Solution [Bhuiyan et al., 2012] *

**Mining step?** Active Preference Learning for Ranking Patterns [Dzyuba et al., 2013] **

*Pattern sampling
**Optimal pattern mining via beam search

Optimal pattern mining

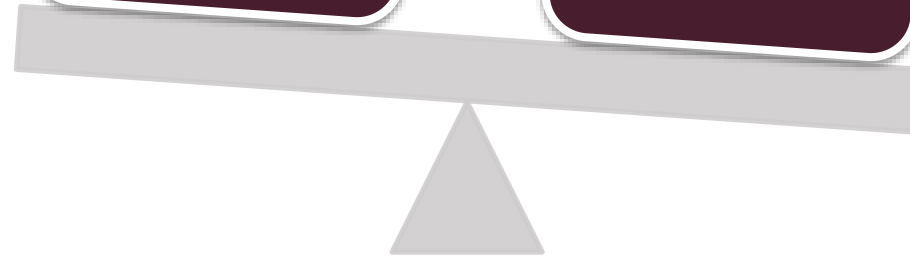Pattern sampling

Very fast

Best patterns

Diversity

Mining step?

Discovering Interesting Patterns Through User's Interactive Feedback [Xin et al., 2006]*

Interactive Pattern Mining on Hidden Data: A Sampling-based Solution [Bhuiyan et al., 2012]**

Active Preference Learning for Ranking Patterns [Dzyuba et al., 2013]*

Learning step?

*Ranking over all patterns = learning to rank problem

**Weight on items

# Conclusion

Frequent pattern mining
1990s

Constraint-based pattern mining
2000s

Optimal pattern mining
Early 2010s

Declarative pattern mining
Early 2010s

Interactive pattern mining
Now

Retrieval era

Exploratory analysis era

Performance issue

The more, the better

Data-driven

Quality issue*

The less, the better*

User-driven

*Faster, better

# Conclusion

Frequent pattern mining
1990s

Constraint-based
pattern mining
2000s

Optimal pattern
mining
Early 2010s

Declarative pattern
mining
Early 2010s

Interactive pattern
mining
Now

Retrieval era

Exploratory analysis era

Performance issue

Quality issue

The more, the better

The less, the better*

Data-driven

User-driven*

*Faster, better, easier

**Thank you!**