



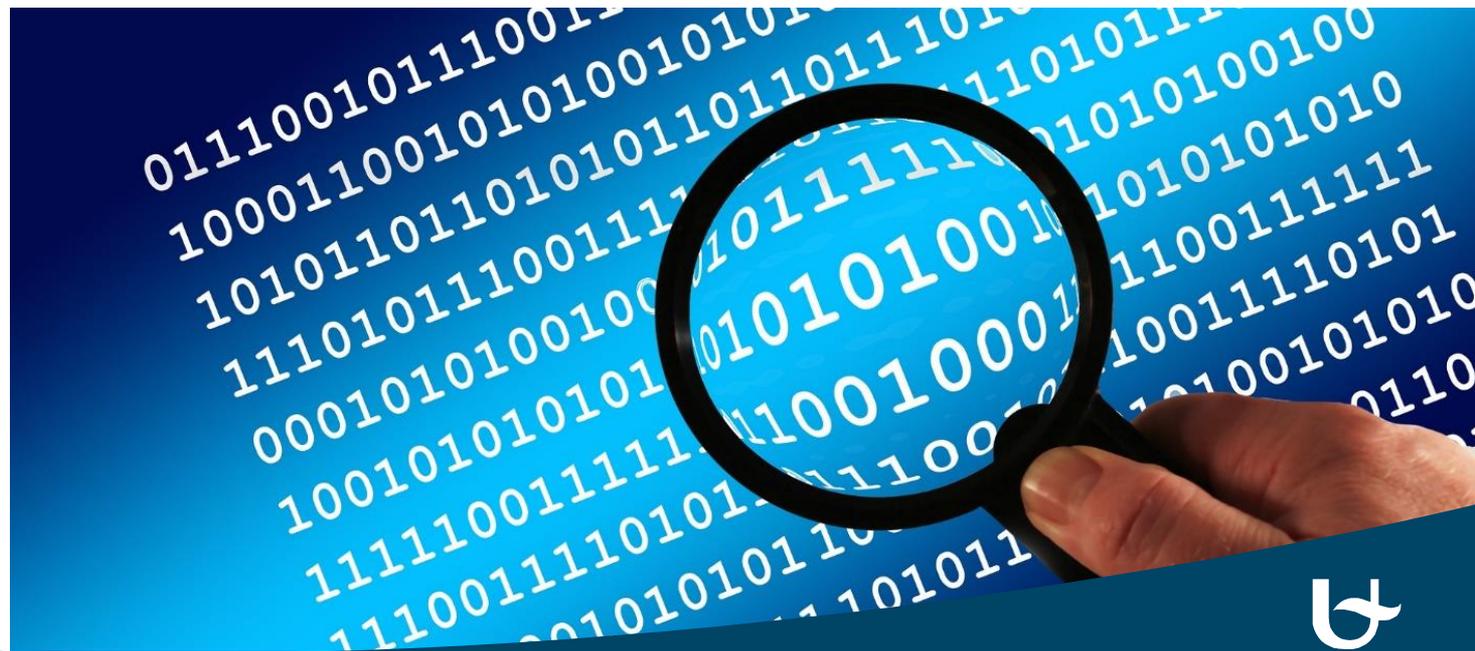
# Data mining, social networks and ethical implications

Toon Calders – Departement Wiskunde-Informatica

# Data Mining

Data mining is the use of automatic techniques to “discover *knowledge*”

- Data driven discovery
- Making implicit knowledge explicit



# Example: Research papers

Can we automatically detect topics one or more researchers are writing about?

Are there remarkable changes in the topics they write about?



# Example: Research papers

## Can we automatically detect topics one or more researchers are writing about?

Efficient Pattern Mining of Uncertain Data with Sampling

Toon Calders<sup>1</sup>, Calin Garboni<sup>2</sup>, and Bart Goethals<sup>2</sup>

<sup>1</sup> TU Eindhoven, The Netherlands

<sup>2</sup> University of Antwerp, Belgium

**Abstract.** Mining frequent itemsets from transactional data is a well known problem with good algorithmic solutions. In the case of uncertain data, however, several new techniques have been proposed. Unfortunately, these proposals often suffer when a lot of items occur with many different probabilities. Here we propose an approach based on sampling by instantiating “possible worlds” of the uncertain data, on which we subsequently run optimized frequent itemset mining algorithms. In this way we gain efficiency at a surprisingly low loss in accuracy. This is confirmed by a statistical and an empirical evaluation on real and synthetic data.

### Abstract

Recent studies on frequent itemset mining algorithms have shown significant performance improvements. However, if the support threshold is set too low, or the data is highly correlated, the number of frequent itemsets itself can be prohibitively large. In this paper, recently several proposals have been made to reduce the size of the frequent itemsets, instead of mining all of them. The main goal of this paper is to identify a subset of all frequent itemsets and to exploit these to reduce the result of a mining operation. We present a new algorithm that derives tight bounds on the support of candidate itemsets. The deduction rules allow for constructing a minimal set of frequent itemsets. We also present connections between recent proposals for concise representations and our experiments on real-life datasets that show the effectiveness of our approach. In fact, the experiments even show that in many cases the concise representation, and then creating the frequent itemsets from this representation outperforms existing frequent itemset mining algorithms.

**Keywords** Discrimination-aware classification · Naive Bayes · Expectation maximization

Computational Complexity of Itemset Mining and Satisfiability

Toon Calders<sup>1</sup>  
toon.calders@ua.ac.be  
Tel: +32 3 265 38 61 Fax: +32 3 265 38 62  
University of Antwerp, Belgium  
Middelheimlaan 1, BE-2020 Antwerp

### ABSTRACT

Computing frequent itemsets is one of the most prominent problems in data mining. We introduce a new, related problem, called **FREQSAT**: given some itemset-interval pairs, does there exist a database such that for every pair the frequency of the itemset falls in the interval? It is shown in this paper that **FREQSAT** is not finitely axiomatizable and that it is **NP**-complete. We also study cases in which other characteristics of the database are given as well. These characteristics can complicate **FREQSAT** even more. For example, when the maximal number of duplicates of a transaction is known, **FREQSAT** becomes **PP**-hard. We describe applications of **FREQSAT** in frequent itemset mining algorithms and privacy in data mining.

following frequency

$freq(a, b)$

There are many itemsets. Exclusion principle

The relationship between the frequency of itemsets and the complexity of the problem comes from the fact that the number of frequent itemsets is exponential in the size of the database. **FREQSAT**: given a set of frequent itemsets, does there exist a database such that the frequency of each itemset is in the given interval?

A Survey

Mining Compressing Sequential Patterns

Hoang Thanh Lam<sup>1</sup>, Fabian Mörchen<sup>2</sup>, Dmitry Fradkin<sup>3</sup> and Toon Calders<sup>1</sup>

<sup>1</sup>Siemens Corporation, Corporate Research and Technology, Princeton, NJ, USA

<sup>2</sup>Amazon.com Inc., 410 Terry Avenue North, Seattle WA 98109, USA

<sup>3</sup>Technische Universiteit Eindhoven, Department of Maths and Computer Science, Eindhoven, Netherlands

Received 29 July 2012; revised 27 February 2013; accepted 16 April 2013

DOI:10.1002/sam.11192

Published online 23 May 2013 in Wiley Online Library (wileyonlinelibrary.com).

**Abstract:** Pattern mining based on data compression has been successfully applied in many data mining tasks. For itemset data, the Krimp algorithm based on the *minimum description length* (MDL) principle was shown to be very effective in solving the redundancy issue in descriptive pattern mining. However, for sequence data, the redundancy issue of the set of frequent sequential patterns is not fully addressed in the literature. In this article, we study MDL-based algorithms for mining non-redundant sets of sequential patterns from a sequence database. First, we propose an encoding scheme for compressing sequence data with sequential patterns. Second, we formulate the problem of mining the most compressing sequential patterns from a sequence database. We show that this problem is intractable and belongs to the class of inapproximable problems. Therefore, we propose two heuristic algorithms. The first of these uses a two-phase approach similar to Krimp for itemset data. To overcome performance issues in candidate generation, we also propose GoKrimp, an algorithm that directly mines compressing patterns by greedily extending a pattern until no additional compression benefit of adding the extension into the dictionary. Since checks for additional compression benefit of an extension are computationally expensive we propose a dependency test which only chooses related events for extending a given pattern. This technique improves the efficiency of the GoKrimp algorithm significantly while it still preserves the quality of the set of patterns. We conduct an empirical study on eight datasets to show the effectiveness of our approach in comparison to the state-of-the-art algorithms in terms of interpretability of the extracted patterns, run time, compression ratio, and classification accuracy using the discovered patterns as features for different classifiers. © 2013 Wiley Periodicals, Inc. *Statistical Analysis and Data Mining* 7: 34–52, 2014

**Keywords:** sequence data; compressing patterns mining; complexity; minimum description length; compression-based pattern mining

this area has been

representations w.r.t. frequency queries. Such representations can be used to support the discovery of every frequent set and its support without looking back at the data. Interestingly, the size of condensed representations can be several orders of magnitude smaller than the size of frequent set collections. Most of the proposals concern exact representations while it is also possible to consider approximated ones, i.e., to trade computational complexity with a bounded approximation on the computed support values. This paper surveys the core concepts used in the recent works on condensed representation for frequent sets.



# Example: Research papers

Can we automatically detect topics researchers are writing about?

Harder than “just counting” frequent subsequences.



Pattern	Support	Pattern	Support
<u>algorithm algorithm</u>	0.376	method method	0.250
<b>learn learn</b>	0.362	<u>algorithm result</u>	0.247
<b>learn algorithm</b>	0.356	Data set	0.244
<u>algorithm learn</u>	0.288	<b>learn learn learn</b>	0.241
data data	0.284	<b>learn</b> problem	0.239
<b>learn</b> data	0.263	<b>learn</b> method	0.229
model model	0.260	<u>algorithm data</u>	0.229
problem problem	0.258	<b>learn</b> set	0.228
<b>learn</b> result	0.255	problem <b>learn</b>	0.227
problem <u>algorithm</u>	0.251	<u>algorithm algorithm algorithm</u>	0.222

Figure 1: The 20 most frequent non-singleton closed sequential patterns from the JMLR abstracts datasets. This set, despite containing some meaningful patterns, is very redundant.

# Example: Research papers

Can we automatically detect topics researchers are writing about?

Harder than “just counting” frequent subsequences.



## Minimal description length principle

A good model learns us a lot about the data and allows to represent it compactly. We hence measure model quality as:

$$\text{Length}(M) + \text{Length}(D | M)$$

# Better Model = More Compression

Biased die



1/2

0



1/4

10



1/8

110



1/16

1110



1/32

11110



1/32

11111

Expected code length:

$$\frac{1}{2} + \frac{1}{4} + \frac{3}{8} + \frac{4}{16} + \frac{5}{32} + \frac{5}{32} = 1.84 \dots$$

Versus without information:

$$\log(6) = 2.58 \dots$$

# Minimal Description Length



Pattern	Support	Pattern	Support
algorithm algorithm	0.376	method method	0.250
learn learn	0.362	algorithm result	0.247
learn algorithm	0.356	Data set	0.244
algorithm learn	0.288	learn learn learn	0.241
data data	0.284	learn problem	0.239
learn data	0.263	learn method	0.229
model model	0.260	algorithm data	0.229
problem problem	0.258	learn set	0.228
learn result	0.255	problem learn	0.227
problem algorithm	0.251	algorithm algorithm algorithm	0.222

$$\sum_{w \in D} \left( \log \frac{F_C}{f_C(w)} * f_C(w) + g_C(w) \right)$$

# Example: Research papers

Can we automatically detect topics researchers are writing about?

Harder than “just counting” frequent subsequences.



Based on Minimal Description Length principle:

support vector machin	larg scale	featur select	sampl size
machin learn	nearest neighbor	graphic model	learn algorithm
state art	decis tree	real world	princip compon analysi
data set	neural network	high dimension	logist regress
bayesian network	cross valid	mutual inform	model select

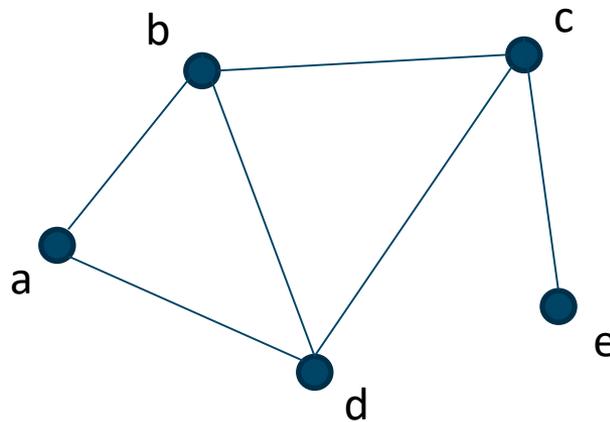




# Neighborhood Function

Count the number of pairs of nodes at distance 1, 2, 3, ...

1: 6  
2: 3  
3: 1



source	dest
a	b
b	c
c	e
d	a
d	c
b	d

Allows to compute average degree, diameter, effective diameter.

# Neighborhood Function

Straightforward algorithm:

Input:  $G(V,E)$

- Let  $N_i(v)$  be the nodes at distance  $i$  of  $v$ ;
- Exploit  $N_{i+1}(v) = \bigcup_{w:\{v,w\}\in E} N_i(w)$

Time:  $O(\delta |E| |V|)$

Space:  $O(|V|^2)$

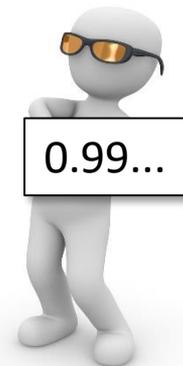


# Tool: Sketch

(Main idea behind Flajolet-Martin sketches)



How many people  
watch my presentation?



Do not keep track of all people, only remember the lowest number

# Tool: Sketch

## (Main idea behind Flajolet-Martin sketches)

How low will this lowest number be?

- The more people there are, the lower we expect this number to be ( $1/(n+1)$  to be exact)
- We can use this number to estimate the number of people that saw the presentation:  $1/\text{min} - 1$

What about the variance?

- Give every person a tuple of  $k$  numbers
- Keep track of  $\min(\text{num}_1), \min(\text{num}_2), \dots, \min(\text{num}_k)$
- Take median of the estimations  $\rightarrow$  large error only if there is a large error in 50% of the estimates

# Neighborhood Function

Based on the sketches of two sets, it is straightforward to create a sketch of their union

Using sketches for  $|N_i(v)|$  instead of  $N_i(v)$  itself brings down complexity to time  $O(\delta |E|)$  and space  $O(|V|)$

- Next to sketching, also sampling, distributed computing
- Solving several social network questions
  - Who influences whom (retweets, repost messages, repeat actions)
  - How does information spread through the network

# Ethical Considerations

“Men between the ages of 16 and 25 are much more likely to be involved in accidents, or be cited for traffic violations.”

**Sam Belden, Insurance.com VP**



*“Taking the gender of the insured individual into account as a risk factor in insurance contracts constitutes discrimination.”*

**March 1, 2011 EU Court of Justice**

**Research Question:** Can we still build profiles taking legal restrictions into account?

# Conclusion

Data Mining for finding patterns in large data collections

Several challenges:

- Finding relevant and orthogonal patterns
  - Minimal Description Length principle
- Data size; e.g., social network analysis
  - Sketching, sampling, approximation
  - Distributed computing
- Ethical aspects
  - What if our classifiers are not “politically correct”.  
Is this our problem?

