

Pattern Sets

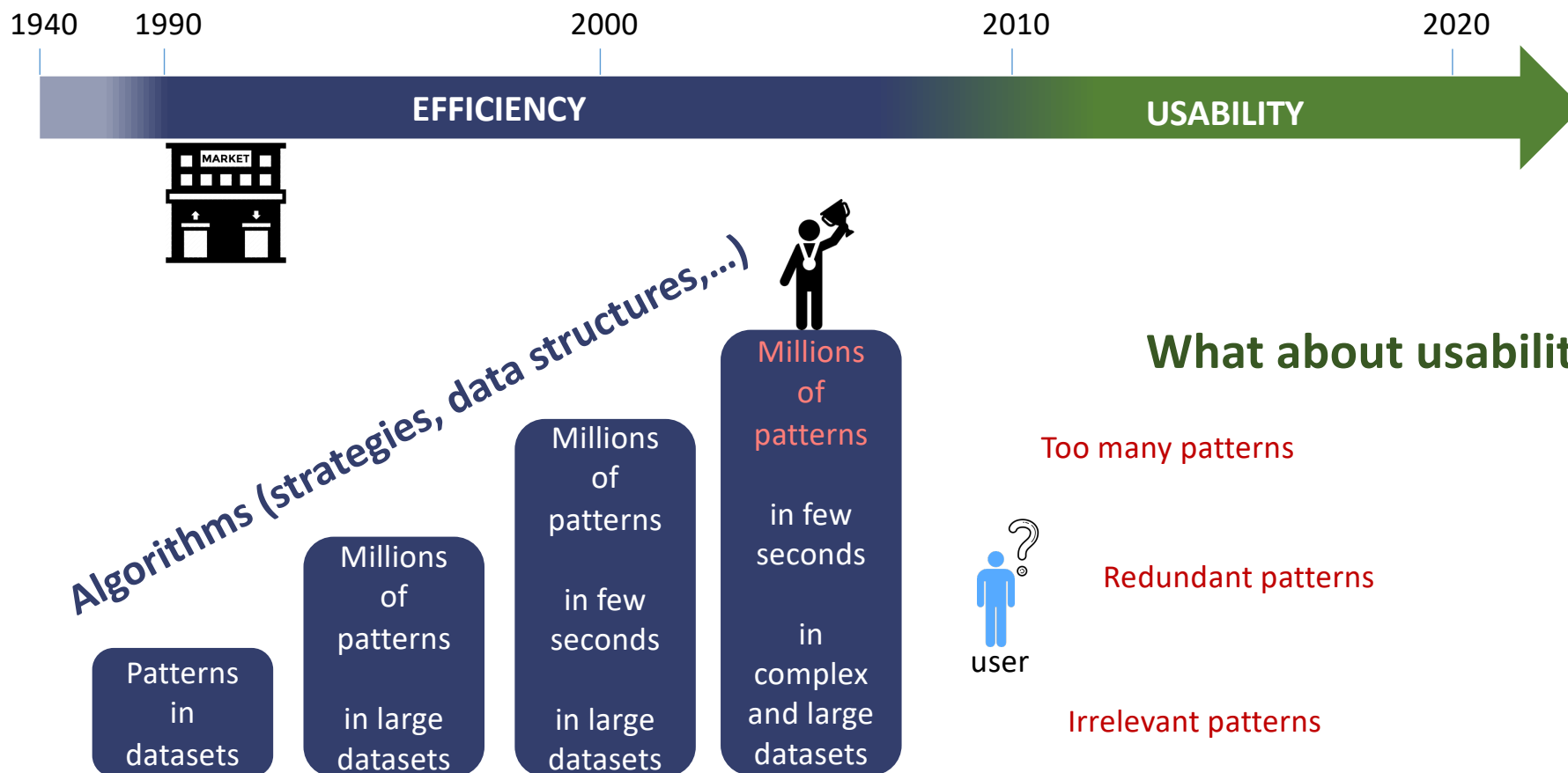
(DMV Lecture, M2 SIF)

Peggy Cellier
peggy.cellier@irisa.fr

Last revision: October 2023

Pattern Mining: Efficiency, Efficiency, Efficiency

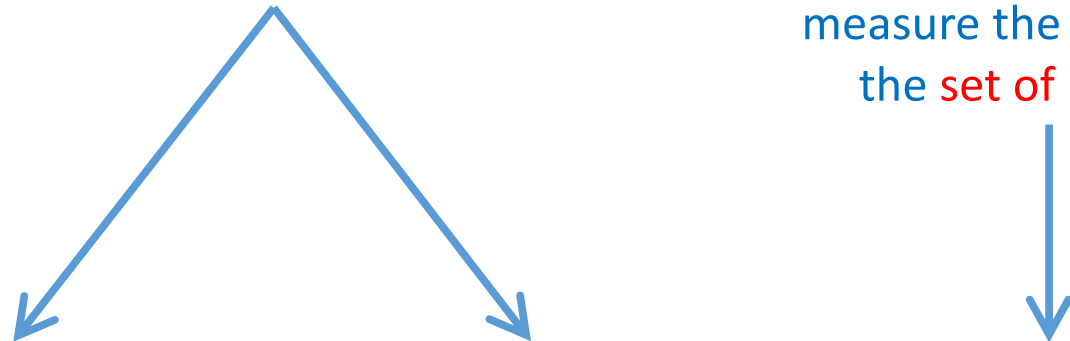
2



How to manage the amount of Patterns?

Most of
measure the interest of
each pattern **individually**

measure the interest of
the **set of patterns**



Condensed representations

Examples

- Closed patterns [Pasquier+'99]
- Maximal patterns

Constraints

Examples

- Length of the patterns
- Regular expressions [Garofalakis+'99]

Measures of interest

Examples [Geng+'06]

- Support
- Confidence
- Growth rate

Pattern sets

Examples

- Tiles [Geerts+'04]
- **Compression-based** [Vreeken+'11]
- Swap randomization [Lijffijt+'12]

References and course material

- [1] Krimp: mining itemsets that compress, Jilles Vreeken, Matthijs van Leeuwen, Arno Siebes, DMKD 2011
- [2] Slim: Directly Mining Descriptive Patterns, Koen Smets and Jilles Vreeken, SDM 2012
- [3] The Long and the Short of It: Summarising Event Sequences with Serial Episodes, Nikolaj Tatti and Jilles Vreeken, KDD 2012
- [4] [The Minimum Description Length Principle](#), Peter Grünwald, MIT Press, 2007
- Some (lot of) slides come from the talk of Jilles Vreeken
 - <https://people.mmci.uni-saarland.de/~jilles/>
- Slides from the HDR of Peggy Cellier

1. Introduction

2. Compression based approaches

1. Motivations

2. Information theory and the MDL principle

- Information theory
- The MDL principle for pattern mining

3. Algorithms

- Itemsets: Krimp (and SLIM)
 - Sequences: SQS
 - Graphs: Graph-MDL
- Conclusion

How can we find **useful patterns?**

For a database db

- a pattern language L
- and a set of constraints C

the **goal** is to find the set of patterns $P \subseteq L$ such that

- each $p \in P$ satisfies each $c \in C$ on db
- and P is maximal

That is, find **all patterns** that satisfy the constraints

The pattern explosion problem

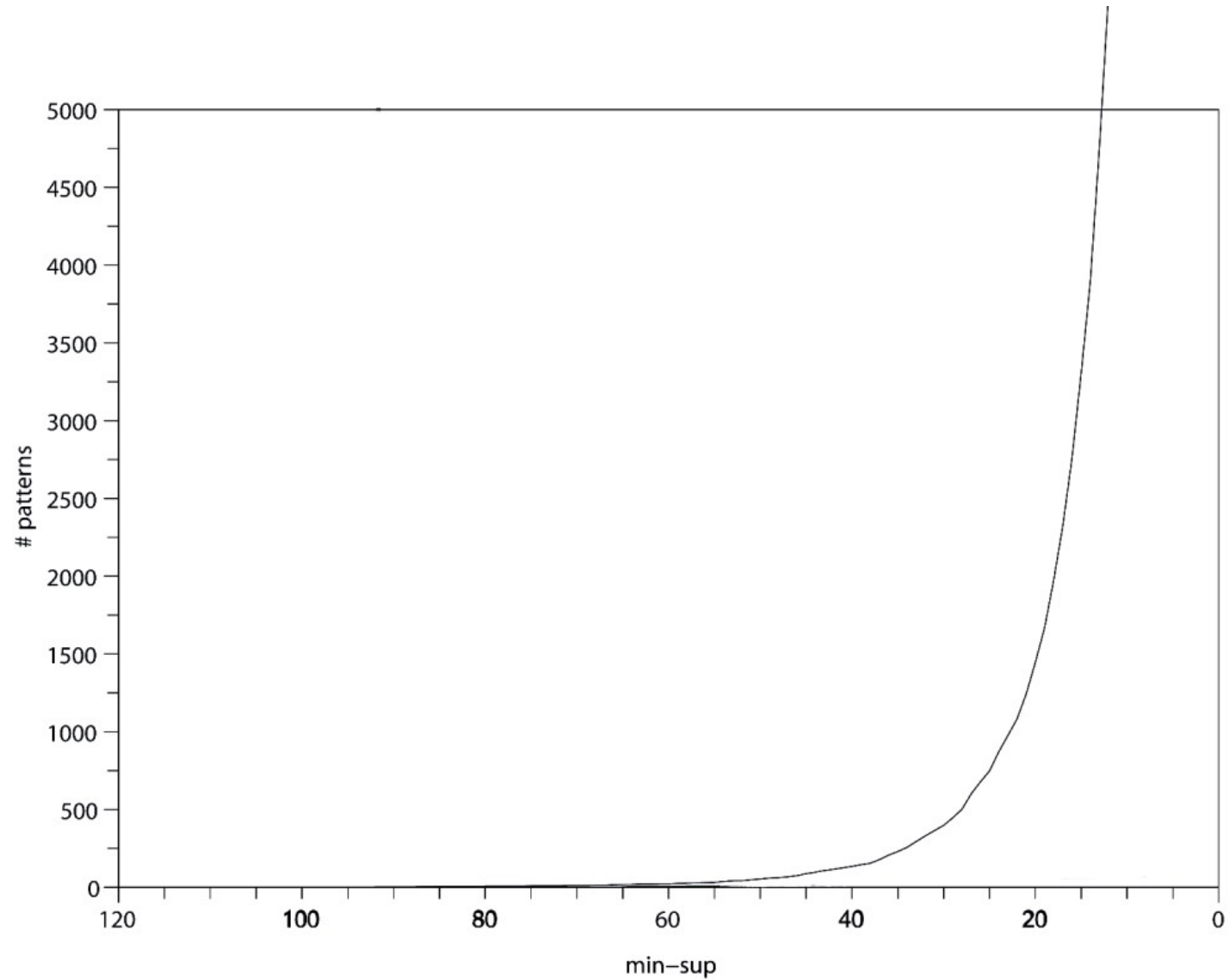
- ↑ high **thresholds**
↓ few, but well-known **patterns**
- ↓ low **thresholds**
↑ a gazillion **patterns**

Many patterns are redundant



Illustration: the *Wine* explosion

- The *Wine* dataset has 178 rows, 14 columns



Be careful what you wish for

The root of all evil is,

- we **ask** for **all** patterns that satisfy some constraints,
- while we want a **small** set that shows the structure of the data

What we ask for

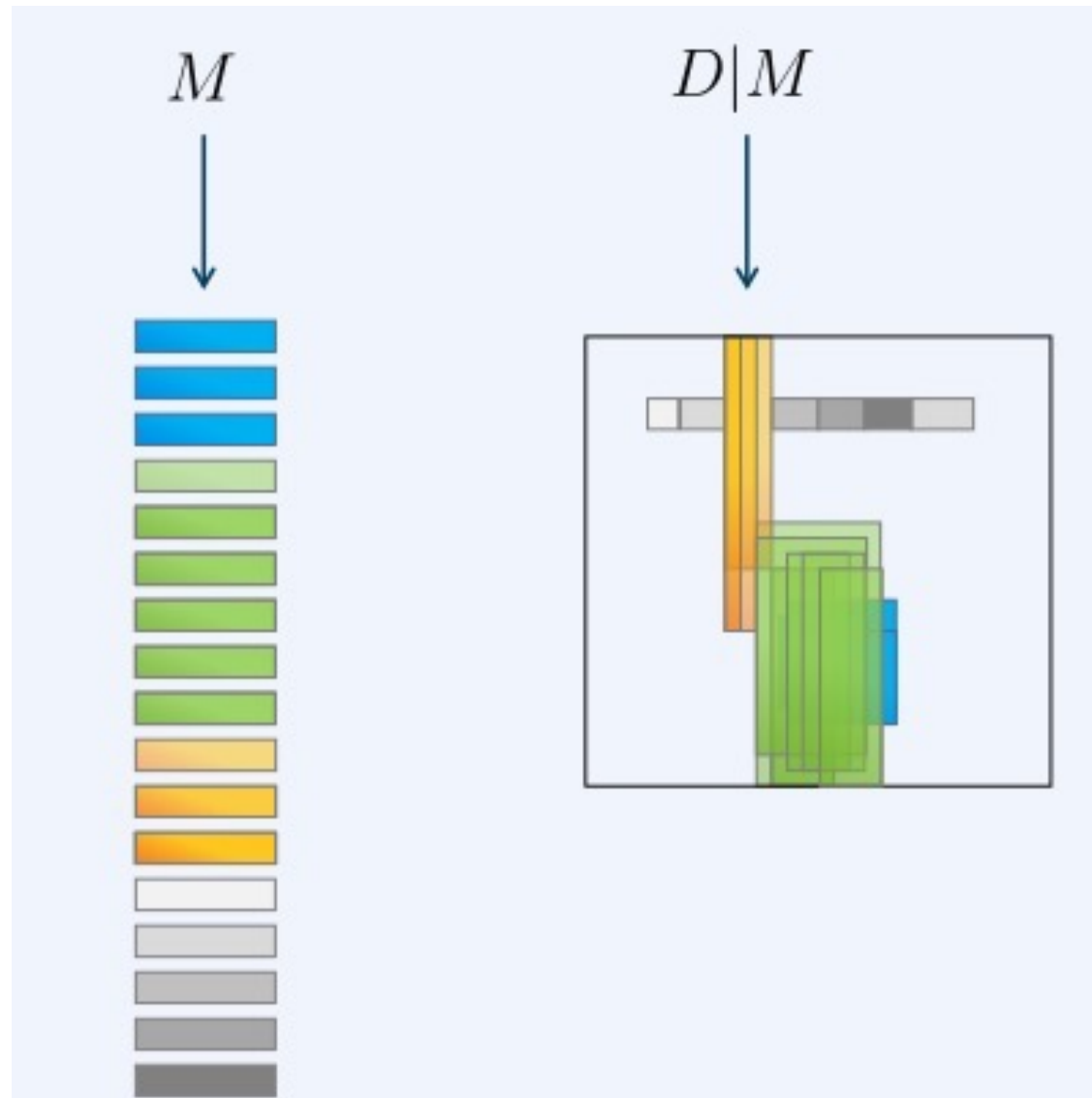
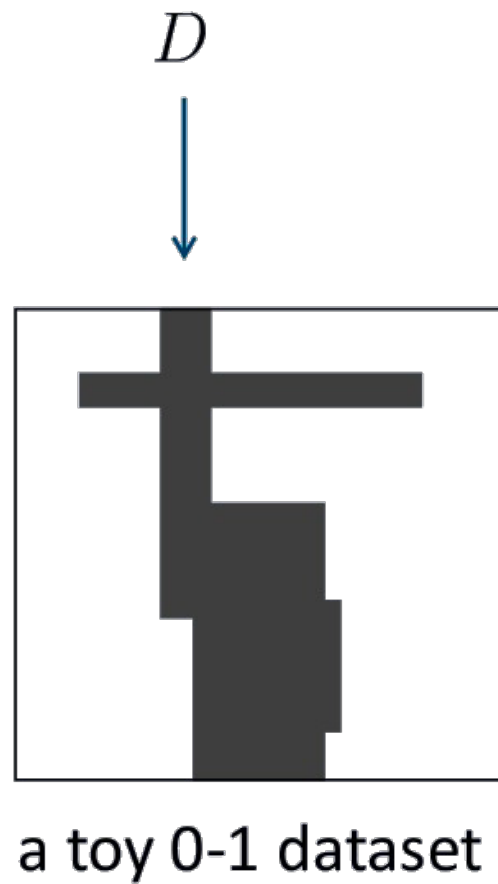
vs

What we really want

In other words, we should ask for a *set of patterns* such that

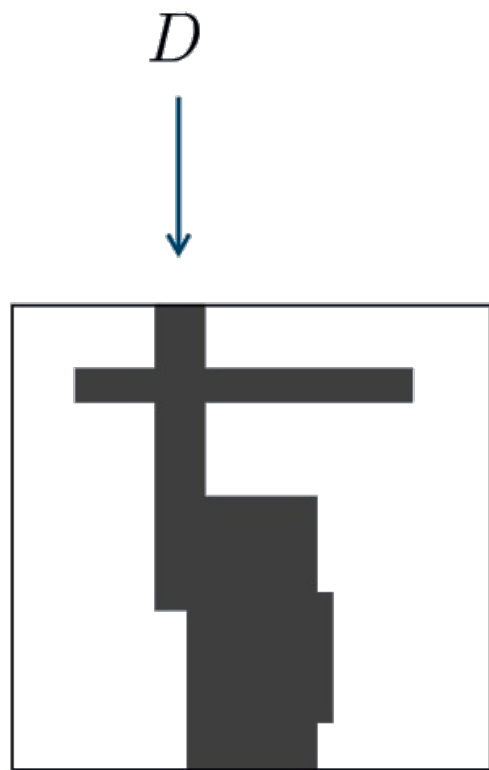
- all members of the set **satisfy** the constraints
- the set is **optimal** with regard to some criterion

What we ask for

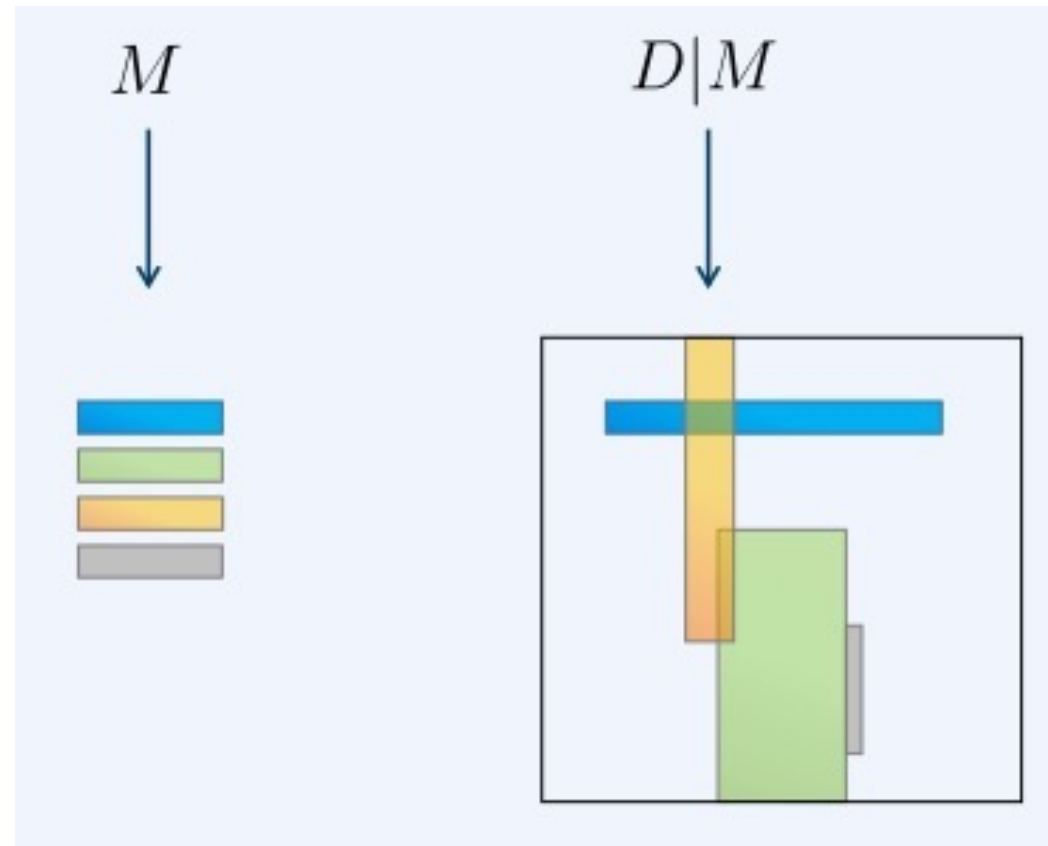


A pattern identifies local properties of the data (e.g., itemsets)

What we really want



a toy 0-1 dataset



How to do that?

How to find

- a subset of patterns
- describing data in a concise way

1. Introduction

2. Compression based approaches

1. Motivations

2. Information theory and the MDL principle

- Information theory
- The MDL principle for pattern mining

3. Algorithms

- Itemsets: Krimp (and SLIM)
- Sequences: SQS
- Graphs: Graph-MDL
- Conclusion

1. Introduction

2. Compression based approaches

1. Motivations

2. Information theory and the MDL principle

- **Information theory**

- The MDL principle for pattern mining

3. Algorithms

- Itemsets: Krimp (and SLIM)
- Sequences: SQS
- Graphs: Graph-MDL

- Conclusion

The Minimum Description Length (MDL) Principle

[Rissanen'78, Grünwald'00]

29

Models describe the data

- that is, they capture regularities
- hence, in an abstract way, they compress it

The MDL principle makes this observation concrete:

“The model that best **describes** the data
is the one that best **compresses** the data.”

The Minimum Description Length (MDL) Principle

[Rissanen'78, Grünwald'00]

30

How to compress data?
⇒ take advantage of repetitions (patterns)

$D = abc c a b c d a b c$



Replace the pattern abc by a short (in number of bits) code ■

$D = \blacksquare c \blacksquare d \blacksquare$

Model (M)

pattern	code
a	■
b	
c	
c	■
d	■

← Pattern appearing **3X** => shorter code

← Patterns appearing **1X** => longer code

$D|M = \blacksquare \blacksquare \blacksquare \blacksquare \blacksquare$





D encoded with M

The Minimum Description Length (MDL) Principle

[Rissanen'78, Grünwald'00]

How to find the BEST compression?
 ⇒ trade-off




Model M'' 1 pattern per item

Pattern	code
a	
b	
c	
d	

D encoded with M''

$D|M'' =$           


Model M

pattern	code
a b c	
c	
d	

D encoded with M

$D|M =$     

Model M' 1 big pattern

pattern	code
a b c c a b c d a b c	

D encoded with M'

$D|M' =$  1 short code

The Minimum Description Length (MDL) Principle

[Rissanen'78, Grünwald'00]

32

How to find the BEST compression?

⇒ trade-off

Principle: Find a trade-off by minimizing

$$L(D,M) = L(M) + L(D|M)$$

Description length of
the model

Description length of
the data encoded with the model

The MDL principle is related to Kolmogorov Complexity

Kolmogorov Complexity [Kolmogorov 1963]

the complexity of a string is the length of the smallest program that generates the string, and then halts

Kolmogorov Complexity is the **ultimate compression**

- recognizes and exploits **any** structure
- Uncomputable however...

1. Introduction

2. Compression based approaches

1. Motivations

2. Information theory and the MDL principle

- Information theory
- **The MDL principle for pattern mining**

3. Algorithms


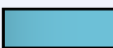



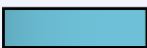
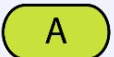

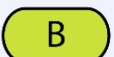

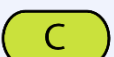




- Itemsets: Krimp (and SLIM)
- Sequences: SQS
- Graphs: Graph-MDL
- Conclusion

- **Goal: reduce the number of extracted patterns**
- **Existing approaches for several pattern languages** [Galbrun'22]
 - Examples
 - Itemsets: KRIMP [Vreeken+'11], SLIM [Smets+'12]
 - Sequences: SQS [Tatti+'12]
 - Graphs: GraphMDL [Bariatti+'20], GraphMDL+ [Bariatti+'21]









Model = Code Table

- **Code Table (CT)**
 - Set of itemsets + encoding

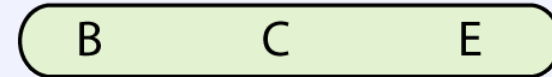
$\mathcal{I} = \{ A, B, C, D, E \}$

Code Table	
<i>Itemset</i>	<i>Code</i>
	
	
	
	
	
	-
	
	





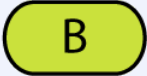



Encoding a database: example

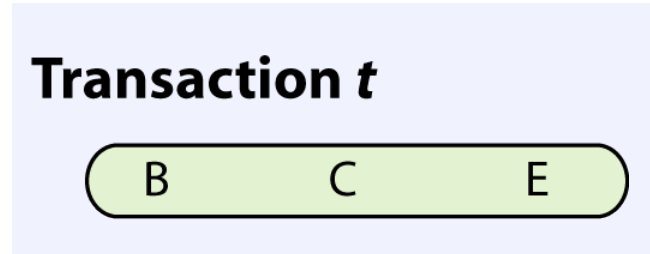
Code Table	
<i>Itemset</i>	<i>Usage</i>
	0
	0
	0
	0
	0
	0
	0
	0

Transaction *t*





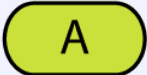
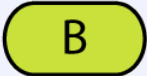





Encoding a database: example

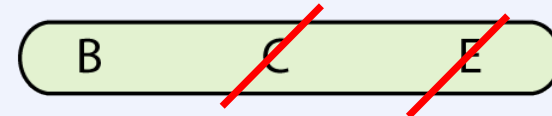
Code Table	
<i>Itemset</i>	<i>Usage</i>
	0
	0
	0
	0
	0
	0
	0
	0



Encoding a database: example

Code Table	
<i>Itemset</i>	<i>Usage</i>
	0
	0
 	0 + 1
	0
	0
	0
	0
	0





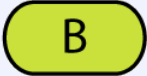



Transaction *t*



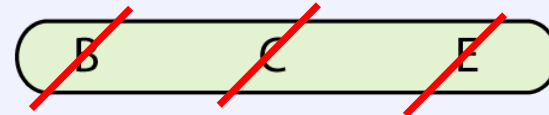
Cover of *t*



Encoding a database: example

Code Table <i>Itemset</i>	<i>Usage</i>
	0
	0
	0 + 1
	0
	0 + 1
	0
	0
	0

Transaction *t*



Cover of *t*



Encoding a database: exercise

- **Exercise:**

- Let us consider the following database and code table
- Find the cover of the transactions
- and the updated usages of the code table

Database

A C E
B
A C
A B C D
B C E
D
B C D E
A B D

Code Table

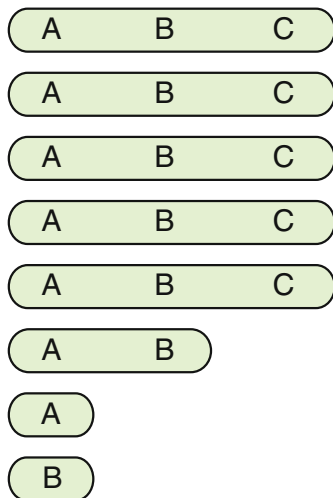
<i>Itemset</i>	<i>Usage</i>
A C	0
B D	0
C E	0
A	0
B	0
C	0
D	0
E	0

Solution

- **ST: Standard Code table**
 - = only and all singleton itemsets

- **Example**

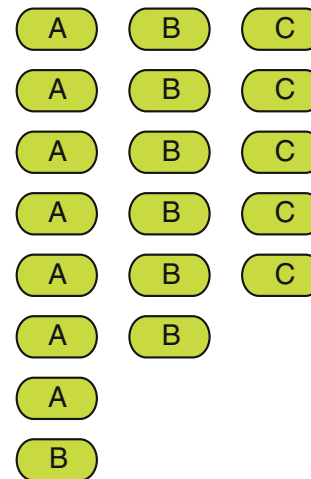
Database



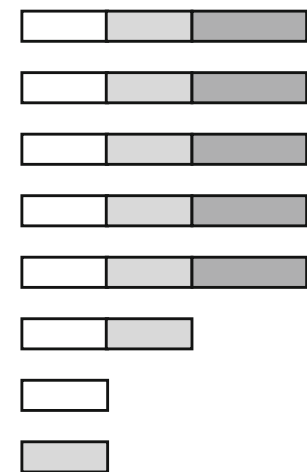
Standard code table *ST*

<i>Itemset</i>	<i>Code</i>	<i>Usage</i>
A		7
B		7
C		5

Cover with *ST*



Encoded with *ST*



- What is the actual length of the codes? (in bits)

- For $c \in CT$ we define

$$P(c|CT) = \frac{usage(c)}{\sum_{i \in CT} usage(i)} \quad \text{where } usage(j) = |\{t \in D \mid j \text{ in } cover(t, CT)\}|$$

- Example: $P(\{A,C\}) = 3/13$

- The optimal code for the coding distribution P assigns a code to $c \in CT$ with length: [Shannon 1948]

$$L(c|CT) = -\log P(c|CT)$$

- Example: $L(\{A,C\}) = -\log(3/13) \approx 2.12$

Length of the code
↔

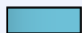



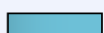


Code Table		Usage	Code
Itemset			
A	C	3	█
B	D	3	█
C	E	2	█
A		1	█
B		2	█
C		0	-
D		1	█
E		1	█

The size of a code table depends on

- the left column
 - length of itemsets as encoded with independence model
 - ST depends only from the data → used to measure CT size
- the right column
 - the optimal code length based on usage of code table elements

Thus, the size of a code table is

$$L(CT | D) = \sum_{c \in CT: usage(c) \neq 0} L(c|ST) + L(c|CT)$$

Code Table		
Itemset	Usage	Code
A C	0	
B D	0	
C E	0	
A	0	
B	0	
C	0	-
D	0	
E	0	

For $t \in D$ we have

$$L(t|CT) = \sum_{c \in \text{cover}(t, CT)} L(c|CT)$$

Description length (DL) for one transaction of D

Hence we have

$$L(D|CT) = \sum_{t \in D} L(t|CT)$$

Sum of the DL of all transactions of D

The total size of data and code table is

$$L(D, CT) = L(CT|D) + L(D|CT)$$

Description length of
the model

Description length of
the data
encoded with the model

- This is the MDL measure that we want to minimize

And now, the optimal code table...

- **Easier said than done**
 - The number of possible code tables is huge
 - No useful structure to exploit

- **Hence, we resort to heuristics**

1. Introduction

2. Compression based approaches

1. Motivations

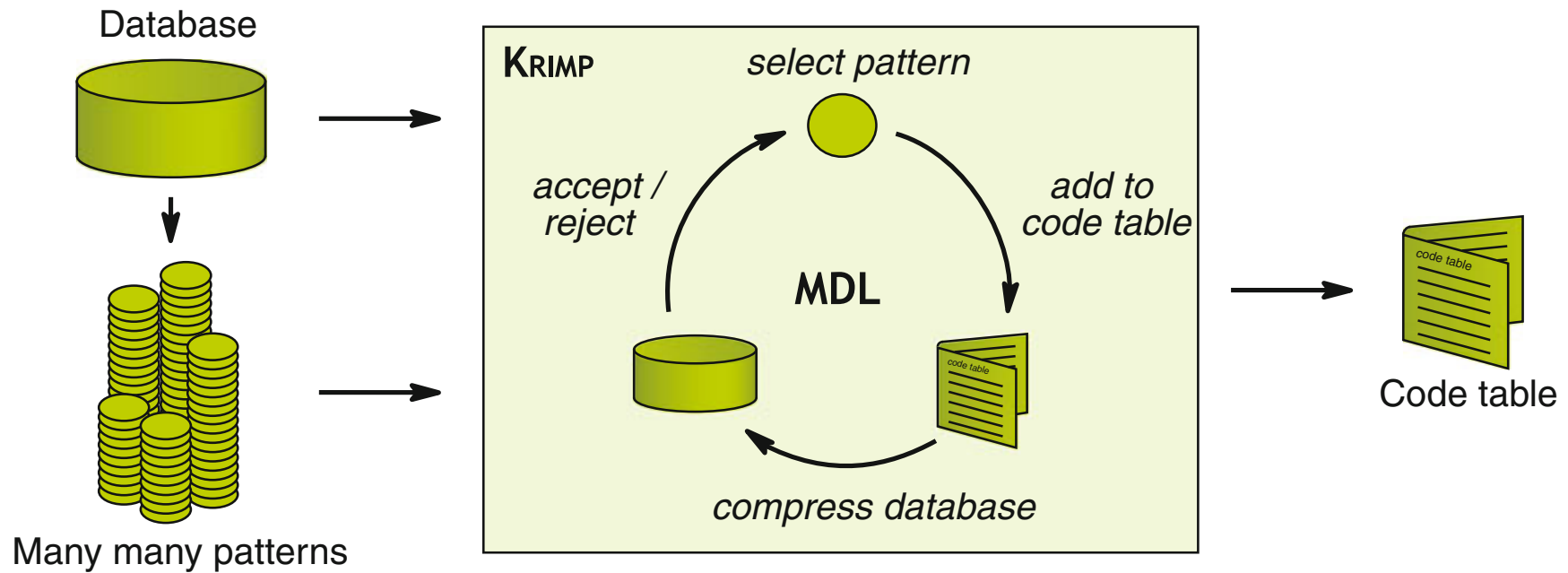
2. Information theory and the MDL principle

- Information theory
- The MDL principle for pattern mining

3. Algorithms

- Itemsets: **Krimp** (and SLIM)
- Sequences: SQS
- Graphs: GraphMDL
- Conclusion

- **KRIMP**
 - Based on MDL
 - And heuristics



Algorithm 3 The KRIMP Algorithm

Input: A transaction database \mathcal{D} and a candidate set \mathcal{F} , both over a set of items \mathcal{I}

Output: A heuristic solution to the Minimal Coding Set Problem, code table CT

- 1: $CT \leftarrow \text{Standard Code Table}(\mathcal{D})$
 - 2: $\mathcal{F}_o \leftarrow \mathcal{F}$ in **Standard Candidate Order** $\text{supp}_{\mathcal{D}}(X) \downarrow$ $|X| \downarrow$ Lexicographically \uparrow
 - 3: **for all** $F \in \mathcal{F}_o \setminus \mathcal{I}$ **do** Add the candidate one by one
 - 4: $CT_c \leftarrow (CT \cup F)$ in **Standard Cover Order** $|X| \downarrow$ $\text{supp}_{\mathcal{D}}(X) \downarrow$ Lexicographically \uparrow
 - 5: **if** $L(\mathcal{D}, CT_c) < L(\mathcal{D}, CT)$ **then**
 - 6: $CT \leftarrow CT_c$
 - 7: **end if**
 - 8: **end for**
 - 9: **return** CT
-

- **Example**

- Let us consider 3 code tables

$$CT_1 = \{\{X_1, X_2\}, \{X_1\}, \{X_2\}, \{X_3\}\}$$

$$CT_2 = \{\{X_1, X_2, X_3\}, \{X_1, X_2\}, \{X_1\}, \{X_2\}, \{X_3\}\}$$

$$CT_3 = \{\{X_1, X_2, X_3\}, \{X_1\}, \{X_2\}, \{X_3\}\}$$

- And assume that $sup(\{X_1, X_2, X_3\}) = sup(\{X_1, X_2\}) - 1$
- KRIMP will never consider CT_3
- But it is possible that $L(D, CT_3) < L(D, CT_2)$
 - => Problem: CT_3 will not be considered

- **Solution: pruning in the code table that KRIMP is considering**

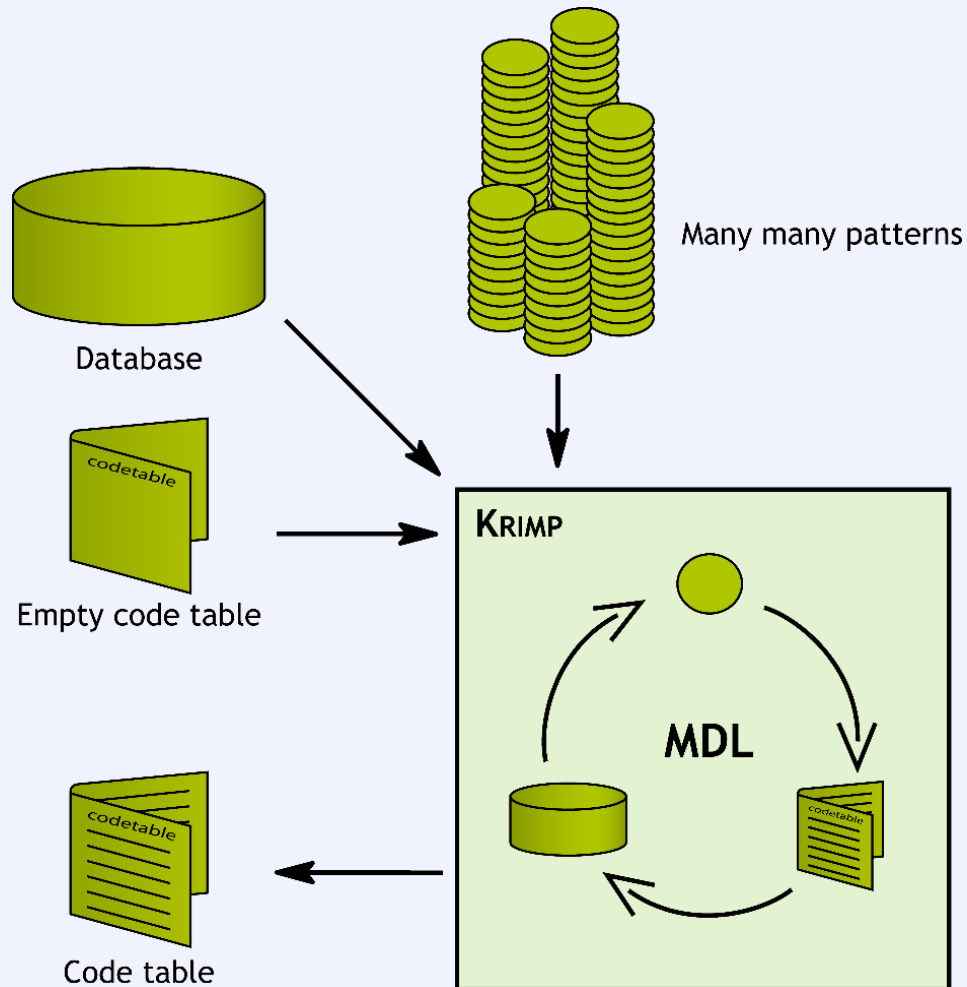
- **When pruning?**
 - When an itemset F is added to the code table
 - At line 6 of KRIMP Algorithm
- **Which itemsets have to be tested?**
 - Only ones for which the usage has decreased
 - Because their code length has increased

Algorithm 4 Code Table Post-Acceptance Pruning

Input: Code tables CT_c and CT , for a transaction database \mathcal{D} over a set of items \mathcal{I} , where $\{X \in CT\} \subset \{Y \in CT_c\}$ and $L(\mathcal{D}, CT_c) < L(\mathcal{D}, CT)$.

Output: Pruned code table CT_p , such that $L(\mathcal{D}, CT_p) \leq L(\mathcal{D}, CT_c)$ and $CT_p \subseteq CT_c$.

- 1: $PruneSet \leftarrow \{X \in CT \mid usage_{CT_c}(X) < usage_{CT}(X)\}$ Only itemsets with decreased usage
 - 2: **while** $PruneSet \neq \emptyset$ **do**
 - 3: $PruneCand \leftarrow X \in PruneSet$ with lowest $usage_{CT_c}(X)$ usage ↓ = code length ↑
 - 4: $PruneSet \leftarrow PruneSet \setminus PruneCand$
 - 5: $CT_p \leftarrow CT_c \setminus PruneCand$
 - 6: **if** $L(\mathcal{D}, CT_p) < L(\mathcal{D}, CT_c)$ **then** Compare the length of \mathcal{D} without the itemset
 - 7: $PruneSet \leftarrow PruneSet \cup \{X \in CT_p \mid usage_{CT_p}(X) < usage_{CT_c}(X)\}$
 - 8: $CT_c \leftarrow CT_p$
 - 9: **end if**
 - 10: **end while**
 - 11: **return** CT_c
-



- mine candidates from D
- iterate over candidates
 - Standard Candidate Order
- covers data greedily
 - no overlap
 - Standard Code Table Order
- select by MDL
 - better compression?
candidates may stay,
reconsider old elements

So, are KRIMP code tables good?

At first glance, yes

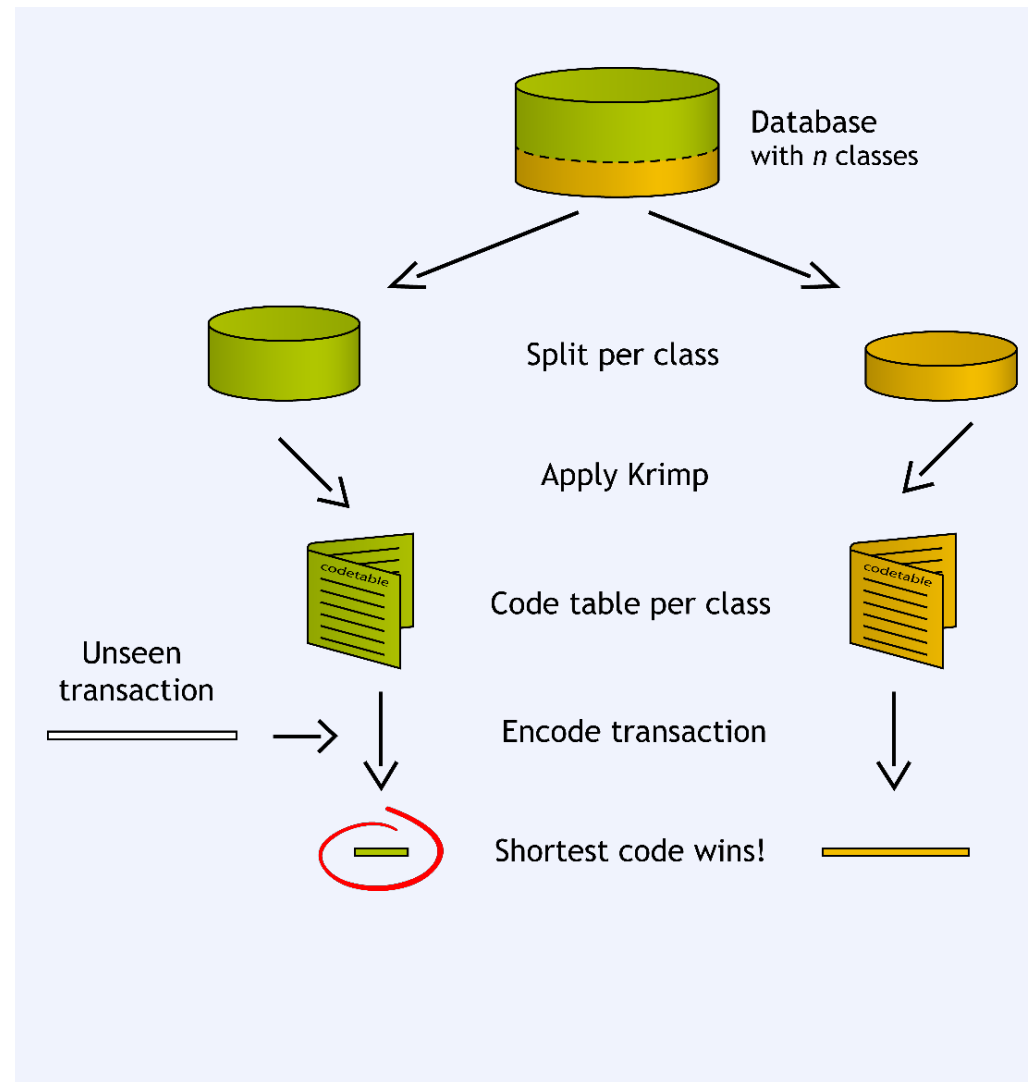
- the code tables are characteristic in the MDL-sense
 - they compress well
- the code tables are small
 - they consist of few patterns
- the code tables are specific
 - they contain relatively long itemsets

But, are these patterns useful?

The quality of the KRIMP code tables was tested by

- classification (ECML PKDD'06)
- measuring dissimilarity (KDD'07, ECML PKDD'15)
- generating data (ICDM'07)
- concept-drift detection (ECML PKDD'08)
- estimating missing values (ICDM'08)
- clustering (ECML PKDD'09)
- sub-space clustering (CIKM'09)
- anomaly detection (SDM'11, CIKM'12, SDM'17)
- characterising uncertain 0-1 data (SDM'11)
- tag-recommendation (IDA'12)
- Web semantic (Semantic Web'20)

- **The Krimp classifier**
 - Split database on class
 - Find code tables
 - Classify by compression



1. Introduction

2. Compression based approaches

1. Motivations

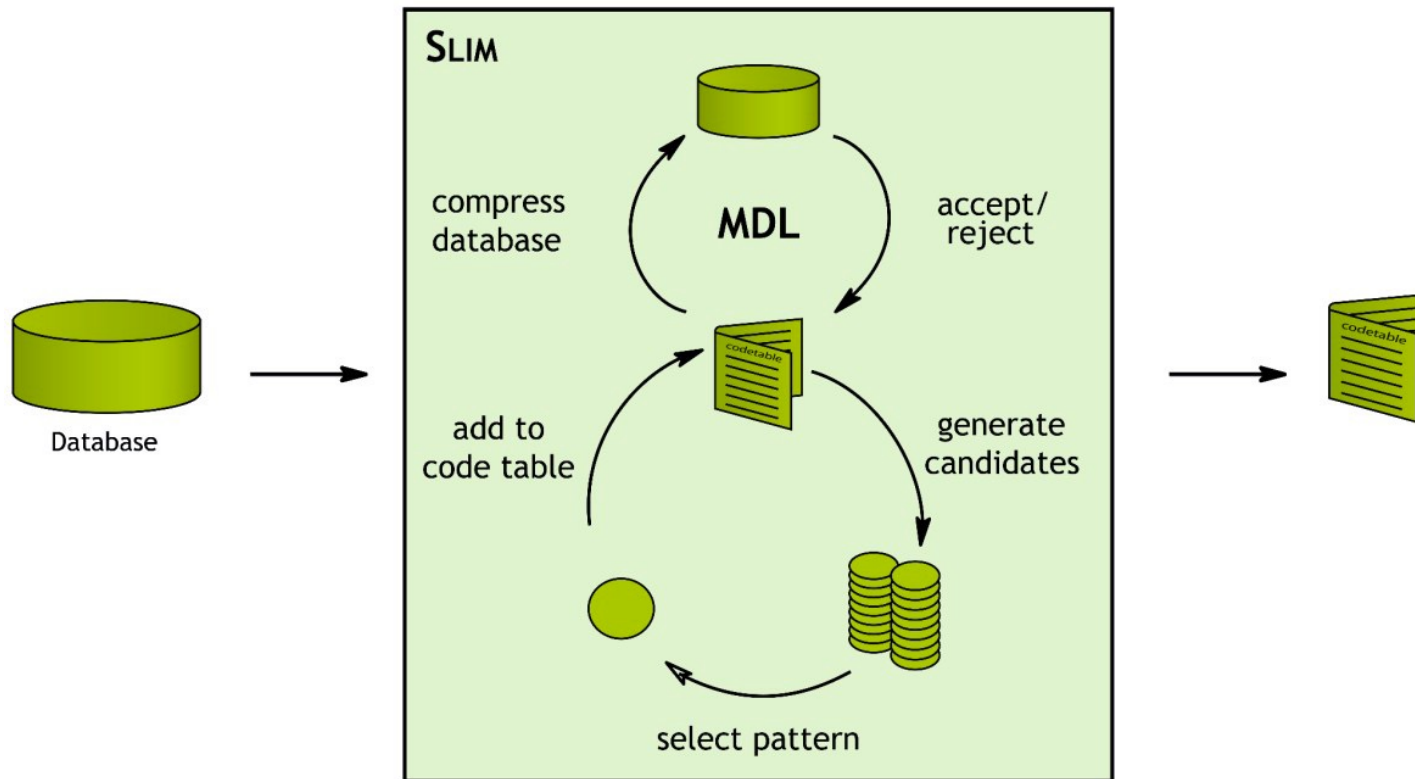
2. Information theory and the MDL principle

- Information theory
- The MDL principle for pattern mining

3. Algorithms

- Itemsets: Krimp (and **SLIM**)
- Sequences: SQS
- Graphs: GraphMDL
- Conclusion

- Any-time version of KRIMP
- Generate-and-select instead of generate-then-select



Algorithm 2 The SLIM Algorithm

Input: A transaction database \mathcal{D} over a set of items \mathcal{I}

Output: A heuristic solution to the Minimal Coding Set Problem, code table CT

1. $CT \leftarrow \mathbf{Standard\ Code\ Table}(\mathcal{D})$
 2. **for** $F \in \{X \cup Y : X, Y \in CT\}$ **in Gain Order** **do** usage↓
 3. $CT_c \leftarrow (CT \oplus F)$ **in Standard Cover Order**
 4. **if** $L(\mathcal{D}, CT_c) < L(\mathcal{D}, CT)$ **then** $|X|$ ↓ $supp_{\mathcal{D}}(X)$ ↓ Lexicographically↑
 5. $CT \leftarrow post-prune(CT_c)$
 6. **end if**
 7. **end for**
 8. **return** CT
-

Note: At line 2 a list of the top-k candidates is kept in case the selected candidate fails the test at line 4

- **Let us assume**
 - A set of items $I=\{A, B, C, D\}$
 - And a database D

transactions	description
T1	A B C
T2	B C D
T3	A B
T4	B C
T5	B D

- **Apply 2 steps of the loop of SLIM on D**
 - For a sake of simplicity
 - No postpruning

Memo:

$$\text{Length}(t) = -\log_2\left(\frac{\text{usage}(t)}{\sum_{y \in CT} \text{usage}(y)}\right)$$

Memo
$-\log_2(5/12) = 1.26$
$-\log_2(3/12) = 2$
$-\log_2(2/12) = 2.58$
$-\log_2(3/9)=1.58$
$-\log_2(2/9) = 2.16$
$-\log_2(1/9) = 3.16$

Solution (to complete)

- 1. Computation of the standard code table

CT = ST =

pattern	support	usage	usage list	length	code

Solution (to complete)

- **4.1 $L(D, CT_c) < L(D, CT)$?**

- $L(D | CT_c) =$

Assume : $L(CT_c | D) = 18.26$

- $L(D | CT) =$

Assume : $L(CT | D) = 15.16$

- **2.2 Selection of a candidate F as the union of 2 elements of CT**

- **3.1 Computation of the new CTc**

Solution (to complete)

53

- **4.1 $L(D, CT_c) < L(D, CT)$?**

- $L(D | CT_c) =$

assume $L(CT_c | D) = 25.74$

- **Characteristics**

- Any-time algorithm
- It considers at each step the refinement of the CT that provides most gain
 - Thanks to heuristics
- Parameter-free (no minsup)

1. Introduction

2. Compression based approaches

1. Motivations

2. Information theory and the MDL principle

- Information theory
- The MDL principle for pattern mining

3. Algorithms

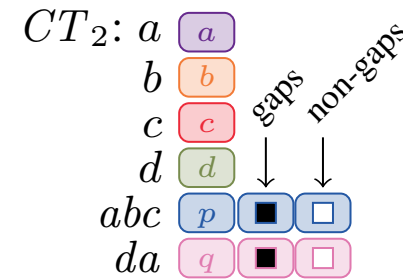
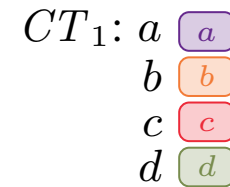
- Itemsets: Krimp (and SLIM)
- Sequences: **SQS**
- Graphs: GraphMDL
- Conclusion

- SQS read « squeeze »
- MDL for sequences

Data D : a b d c a d b a a b c

- **Code table**

- 4 columns
 - Pattern,
 - code,
 - gap code,
 - no gap code



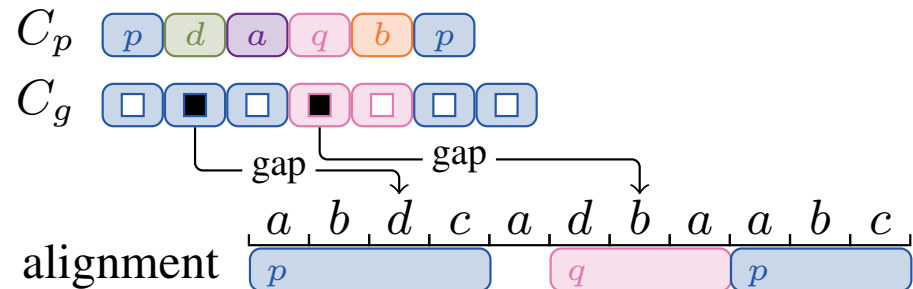
- **Encoded database = 2 streams**

- C_p : pattern stream
- C_g : gap stream

Encoding 1: using only singletons

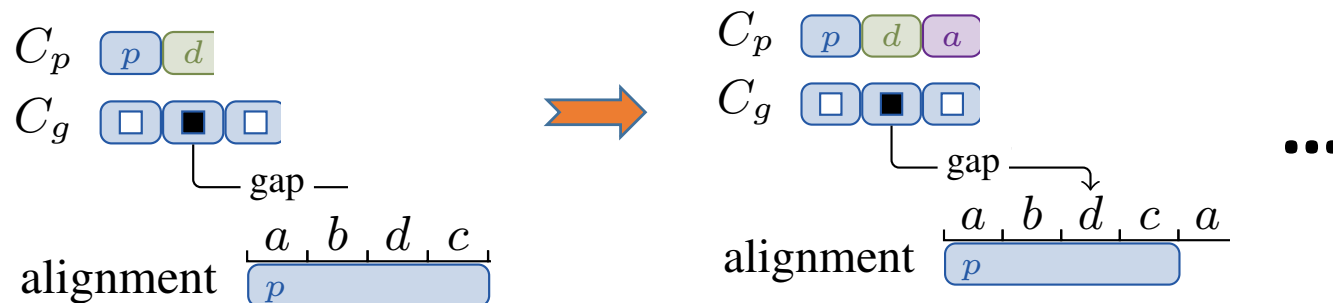
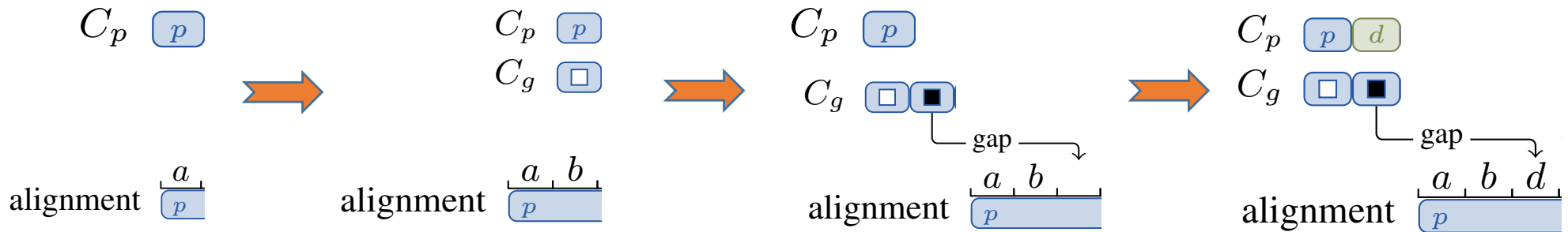
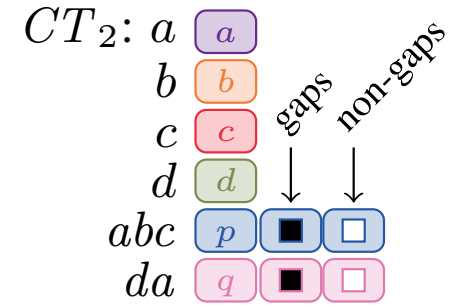
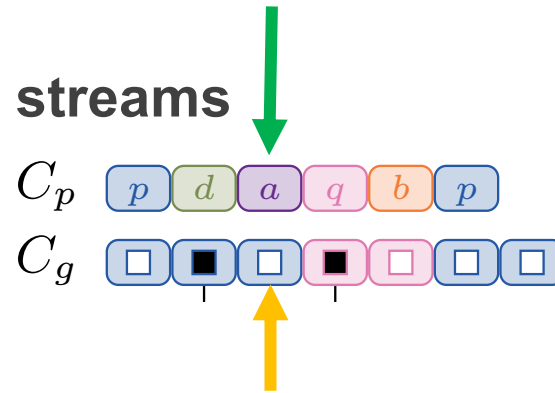


Encoding 2: using patterns

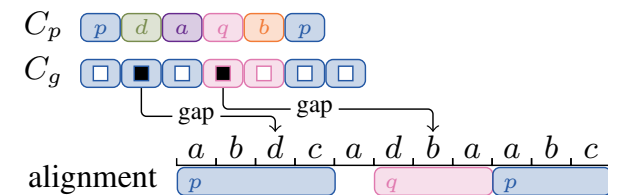


- **Encoded database = 2 streams**

- C_p : pattern stream
- C_g : gap stream



Encoding 2: using patterns



Exercise

- Let us consider the following code table and the following streams
- Give the decoded sequence

CT =

a	a
b	b
c	c
d	d
b c	m ■ □
a b d	n ■ □
c d	o ■ □

Cp	a	m	n	d	b	n	a	d	o
Cg	□	■	□	□	■	■	□	□	□

Encoding length

- Patterns

$$L(\text{code}_p(X) \mid CT) = -\log \left(\frac{\text{usage}(X)}{\sum_{Y \in CT} \text{usage}(Y)} \right)$$

- Let

- $\text{gaps}(X)$ = number of gaps in the usage of pattern X
- $\text{fills}(X)$ = number of non-gaps in the usage of pattern X
 - Note: $\text{fills}(X) = \text{usage}(X)(|X| - 1)$

- Gaps

$$L(\text{code}_g(X) \mid CT) = -\log \left(\frac{\text{gaps}(X)}{\text{gaps}(X) + \text{fills}(X)} \right)$$

- Non-gaps

$$L(\text{code}_n(X) \mid CT) = -\log \left(\frac{\text{fills}(X)}{\text{gaps}(X) + \text{fills}(X)} \right)$$

Encoding length

- Pattern stream

$$L(C_p | CT) = \sum_{X \in CT} usage(X) L(code_p(X))$$

- Gap stream

$$L(C_g | CT) = \sum_{\substack{X \in CT \\ |X| > 1}} \left(gaps(X) L(code_g(X)) + fills(X) L(code_n(X)) \right)$$

- Database D with a code table CT

$$L(D | CT) = L_{\mathbb{N}}(|D|) + \sum_{S \in D} L_{\mathbb{N}}(|S|) + L(C_p | CT) + L(C_g | CT)$$

Number of
sequences in D

Length of a
sequence S in D

Just for information but not to learn: encoding of the code table

- Length of a pattern X

$$L(X, CT) = L_{\mathbb{N}}(|X|) + L_{\mathbb{N}}(\text{gaps}(X) + 1) + \sum_{x \in X} L(\text{code}_p(x \mid ST))$$

Length of the code table

$$L(CT \mid C) = L_{\mathbb{N}}(|\Omega|) + L_U(|D|, |\Omega|) +$$

$$L_{\mathbb{N}}(|\mathcal{P}| + 1) + L_{\mathbb{N}}(\text{usage}(\mathcal{P}) + 1) +$$

$$L_U(\text{usage}(\mathcal{P}), |\mathcal{P}|) + \sum_{X \in \mathcal{P}} L(X, CT)$$

- **Same idea as KRIMP**
 - Input:
 - a set of candidate patterns extracted by an other pattern mining algorithm
 - Method
 - Try to add one by one each candidate and see if the compression is better

Algorithm 3: SQS-CANDIDATES(\mathcal{F}, D)

input : candidate patterns \mathcal{F}

output : set of non-singleton patterns \mathcal{P} that heuristically minimise the **Minimal Code Table Problem**

Sort patterns

With pattern, gain?

- 1 order patterns $X \in \mathcal{F}$ based on $L(D, \{X\})$;
- 2 $\mathcal{P} \leftarrow \emptyset$;
- 3 **foreach** $X \in \mathcal{F}$ in order **do**
- 4 **if** $L(D, \mathcal{P} \cup X) < L(D, \mathcal{P})$ **then**
- 5 | $\mathcal{P} \leftarrow \text{PRUNE}(\mathcal{P} \cup X, D, \text{false})$; Pruning step
- 6 $\mathcal{P} \leftarrow \text{PRUNE}(\mathcal{P}, D, \text{true})$;
- 7 order patterns $X \in G$ by $L(D, \mathcal{P}) - L(D, \mathcal{P} \setminus X)$;
- 8 **return** \mathcal{P} ;

- **Same idea as SLIM**
 - Any-time algorithm
 - Free-parameter

Algorithm 6: SQS-SEARCH(D)

input : database D
output : significant patterns \mathcal{P}

- 1 $\mathcal{P} \leftarrow \emptyset; A \leftarrow \text{SQS}(D, \emptyset);$
- 2 **while** changes **do**
- 3 $\mathcal{F} \leftarrow \emptyset;$
- 4 **foreach** $P \in CT$ **do** add ESTIMATE(P, A, D) to $\mathcal{F};$
- 5 **foreach** $X \in \mathcal{F}$ ordered by the estimate **do**
- 6 **if** $L(D, \mathcal{P} \cup X) < L(D, \mathcal{P})$ **then**
- 7 $\mathcal{P} \leftarrow \text{PRUNE}(\mathcal{P} \cup X, D, \text{false});$
- 8 **if** X is added **then** test recursively X augmented with events occurring in the gaps;
- 9 $\mathcal{P} \leftarrow \text{PRUNE}(\mathcal{P}, D, \text{true});$
- 10 order patterns $X \in G$ by $L(D, \mathcal{P}) - L(D, \mathcal{P} \setminus X);$
- 11 **return** $\mathcal{P};$

1. Introduction

2. Compression based approaches

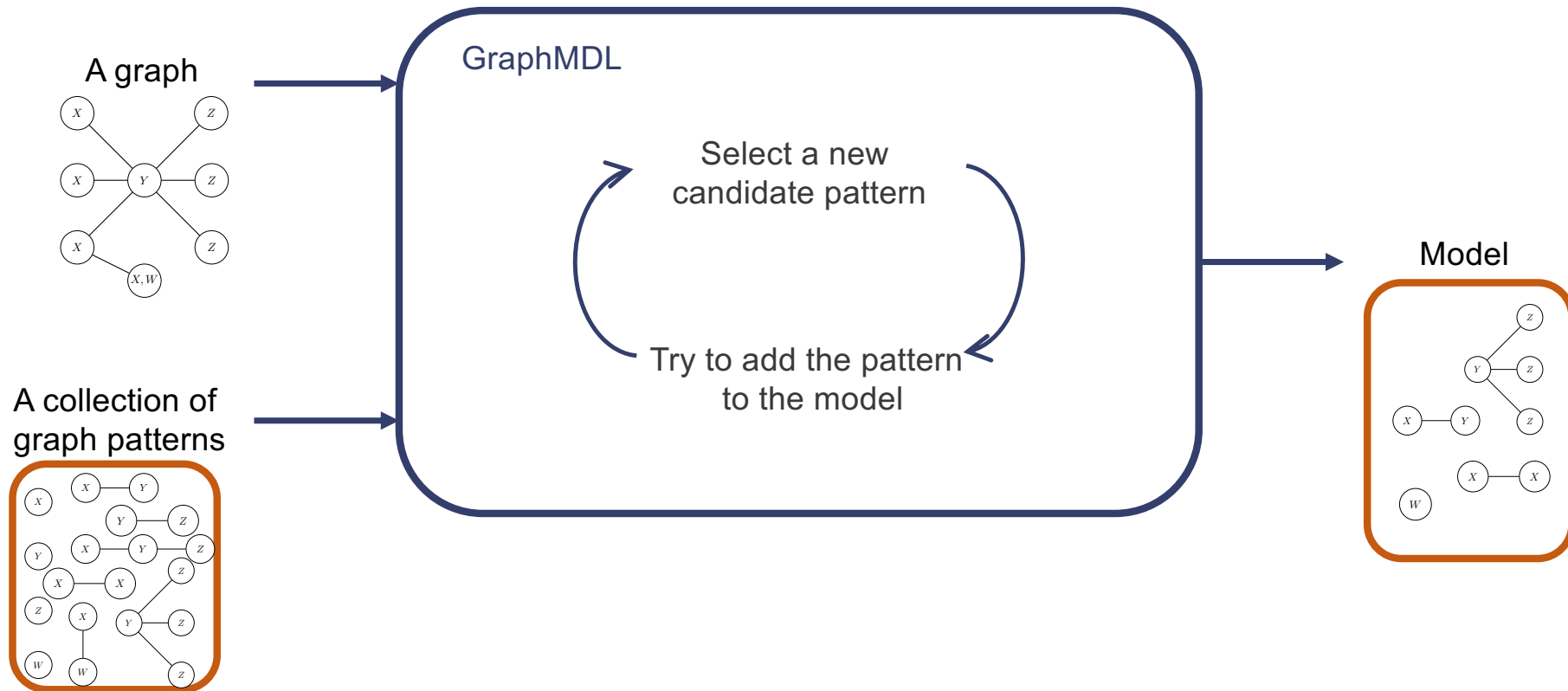
1. Motivations

2. Information theory and the MDL principle

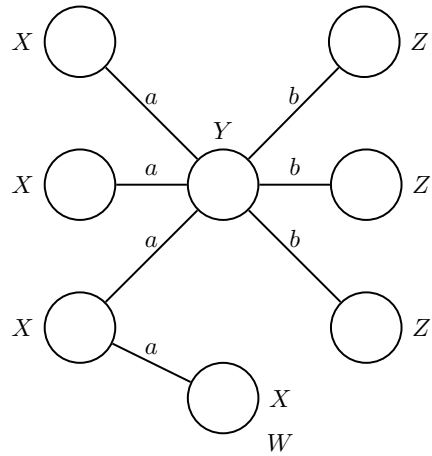
- Information theory
- The MDL principle for pattern mining

3. Algorithms

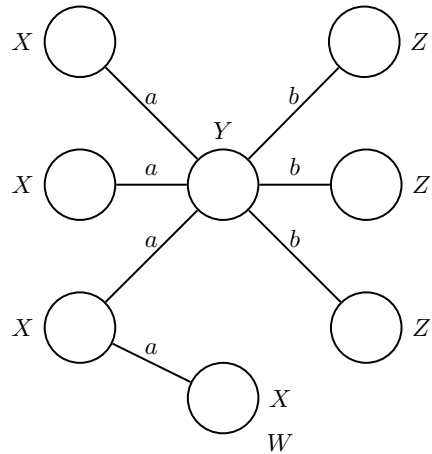
- Itemsets: Krimp (and SLIM)
 - Sequences: SQS
 - Graphs: **GraphMDL**
- Conclusion



- Graph



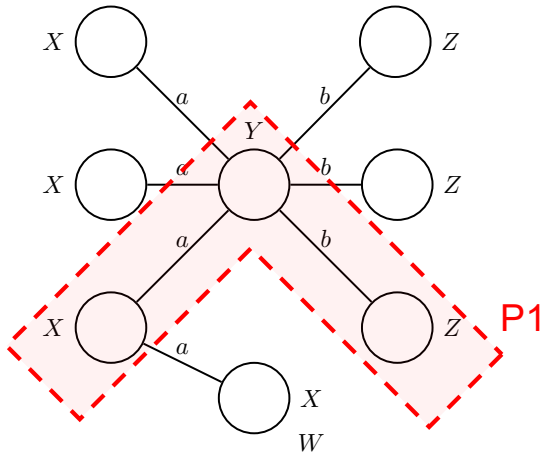
- Graph



Model = Code table

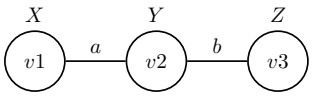

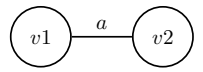

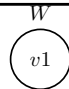

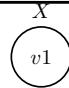

	Pattern structure	Pattern code
P1		P 1
Pa		Pa
Pw		Pw
Px		Px

- Graph

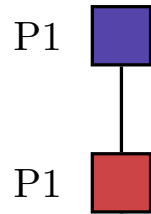
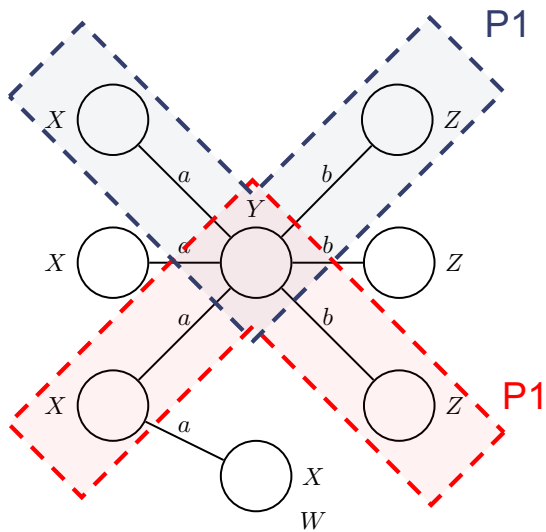


P1 

Model = Code table

	Pattern structure	Pattern code
P1		
Pa		
Pw		
Px		

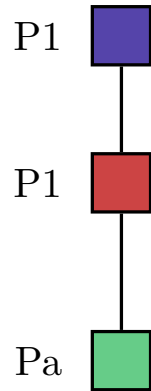
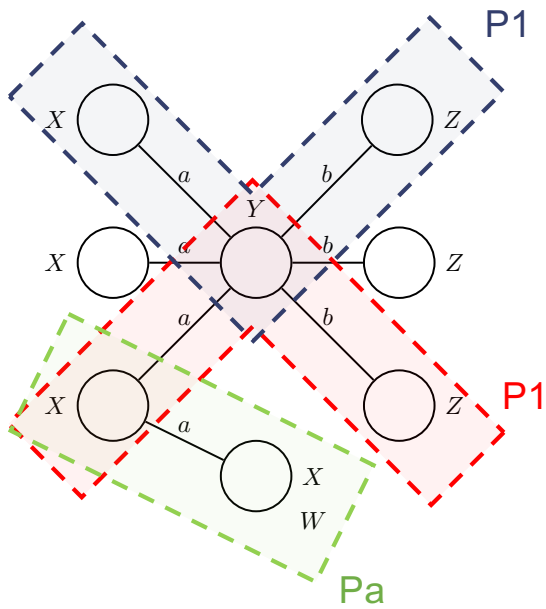
- Graph



Model = Code table

	Pattern structure	Pattern code
$P1$		$P1$
Pa		Pa
Pw		Pw
Px		Px

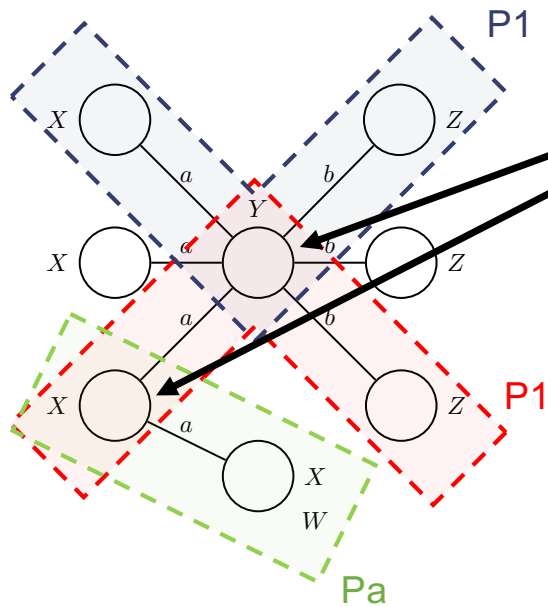
- Graph



Model = Code table

	Pattern structure	Pattern code
$P1$		P_1
Pa		Pa
Pw		Pw
Px		Px

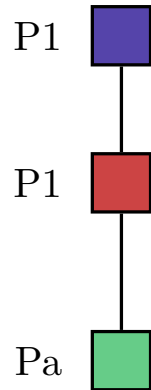
- Graph



Model = Code table

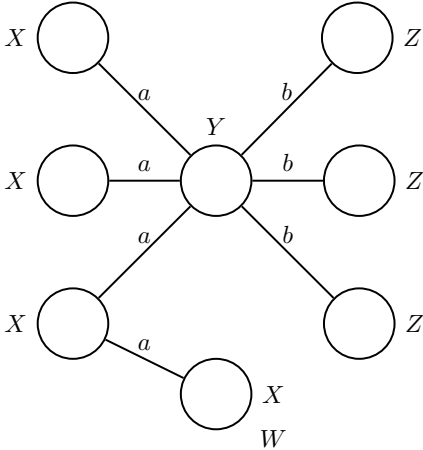
2 different points of connexion

	Structure	Pattern code
P1		
Pa		
Pw		
Px		



BUT no differences in the encoded data!
 ⇒ **Loss of information**

- Graph



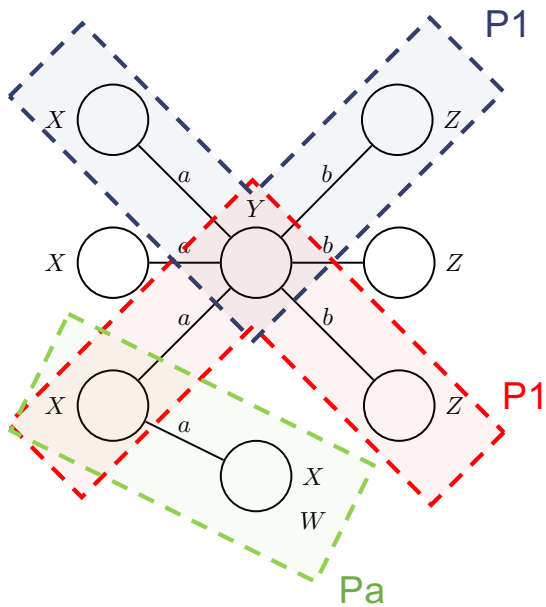
Model = Code table

	Pattern structure	Pattern code	Port ID	Port code
P1			v1 v2	
Pa			v1 v2	
Pw			v1	
Px			v1	

Ports

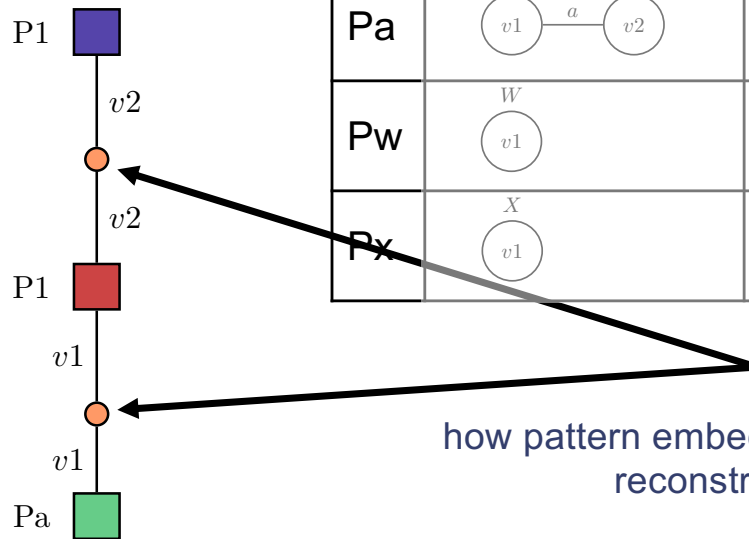
how pattern embeddings connect to each other to reconstruct the original graph

- Graph



Model = Code table

	Pattern structure	Pattern code	Port ID	Port code
P1			v1 v2	
Pa			v1 v2	
Pw			v1	
Px			v1	



Ports

how pattern embeddings connect to each other to reconstruct the original graph

- **3 datasets**

- AIDS-CA, AIDS-CM: about molecules (few labels, many cycles)
 - National Cancer Institute AIDS antiviral seen data
- UD-PUD-En: about dependency trees of english sentences (many labels, no cycles)
 - Universal dependencies project

- **Quantitative results**

Dataset	gSpan minsup	# Candidate patterns input	# Selected patterns output
AIDS-CA	20%	2194	115
AIDS-CA	15%	7867	123
AIDS-CA	10%	20596	148
AIDS-CM	20%	433	111
AIDS-CM	15%	779	131
AIDS-CM	10%	2054	163
AIDS-CM	5%	9943	225
UD-PUD-En	10%	164	162
UD-PUD-En	5%	458	249
UD-PUD-En	1%	6021	523
UD-PUD-En	0%	233434	773

- Drastic reduction
- Relatively stable in the number of selected patterns

- **GraphMDL+ as SLIM and SQS-Search**
 - Generate-and-select algorithm
 - Any-time algorithm
 - Free-parameter

1. Introduction

2. Compression based approaches

1. Motivations

2. Information theory and the MDL principle

- Information theory
- The MDL principle for pattern mining

3. Algorithms

- Itemsets: Krimp (and SLIM)
 - Sequences: SQS
 - Graphs: GraphMDL
- Conclusion

If we have time
Mining Periodic Patterns with a MDL
Criterion

**Esther Galbrun, Peggy Cellier, Alexandre Termier, Bruno
Crémilleux**

- **MDL is great for picking important and useful patterns**
- **KRIMP, SQS, GRAPHMDL approximate the MDL ideal **very well****
 - **vast** reduction of the number of itemsets
 - works for other pattern types equally well: itemsets, sequences, streams
- **Local patterns and information theory**
 - naturally induce good classifiers, clusterers, distance measures
 - with **instant** characterisation and explanation,
 - and, **without** (explicit) parameters

- **When using MDL for pattern mining**
 - Question 1: what kind of model to use?
 - Data modelisation
 - Question 2: how to encode the database? the model?
 - Information theory
 - Question 3: how to find the « best » (or at least one good) model?
 - Algorithmic

1. Introduction

2. Compression based approaches

1. Motivations

2. Information theory and the MDL principle

- Information theory
- The MDL principle for pattern mining

3. Algorithms

- Itemsets: Krimp (and SLIM)
 - Sequences: SQS
 - Graphs: Graph-MDL
- Conclusion