

Lecture 2. White- and black-box models. Causality

Regression models • Decision trees Itemsets mining • Discretization • Clustering • Numerical pattern mining • Tree ensembles • Neural networks

Part I. Pattern mining

(or white-box models in unsupervised settings)

Interpretable models in unsupervised settings

Pattern mining

Goal: to discover a small set of non-redundant and interesting patterns that describe almost entirely the whole dataset

Main steps:

- discover the pattern search space (candidates to interesting patterns)

First of all, we need to select a pattern type. Then, all possible patterns of the chosen type make the pattern search space

- selection of the most interesting ones

We consider pattern mining in **binary** and **numerical** data

Pattern mining in binary data

Transactional dataset and its binary representation

		rice	bread	milk	pasta	yogurt	juice	apples	chips
t ₁	rice, bread, milk	X	X	X					
t ₂	rice, pasta, milk	X		X	X				
t ₃	yogurt, juice, apples, pasta				X	X	X	X	
t ₄	chips, apples, rice	X						X	X
t ₅	pasta, milk, apples			X	X				
t ₆	pasta, juice, bread, yogurt		X		X	X	X		

Itemset: any combination of binary attributes

Formal concepts (an intuitive introduction)

One way to restrict the pattern search space is to use the **maximal rectangles** in data filled with crosses

The rectangle is maximal if we cannot add neither column nor a row such that a new rectangle is filled entirely with the crosses

These maximal rectangles are called **formal concepts**

	rice	bread	milk	pasta	yogurt	juice	apples	chips
t ₁	X	X	X					
t ₂	X		X	X				
t ₃				X	X	X	X	
t ₄	X						X	X
t ₅			X	X				
t ₆		X		X	X	X		

Formal concept analysis

A formal context is a triple (G, M, I) , where G is a set of objects, M is a set of attributes, and I is the incidence relation between two sets, i.e., $(g, m) \in I$ if object g has attribute m

$$A' = \{ m \in M \mid gIm \text{ for all } g \in A \}$$

$$B' = \{ g \in G \mid glm \text{ for all } m \in B \}$$

$$\{t_3, t_6\}' = \{\text{pasta, yogurt, juice}\} \quad A' = B$$

$$\{\text{pasta, yogurt, juice}\}' = \{t_3, t_6\} \quad B' = A$$

Formal concept is a pair (A, B) where

$A' = B$, and $B' = A$. B is called **closed itemset**

	rice	bread	milk	pasta	yogurt	juice	apples	chips
t_1	x	x	x					
t_2	x		x	x				
t_3				x	x	x	x	
t_4	x						x	x
t_5			x	x				
t_6		x		x	x	x		

Specificity: for noisy data the number of patterns may grows exponentially

Formal concept analysis. Closed itemsets

For an arbitrary set of attributes C the formal concept is computed as follows:

$$C' = B$$

$$B' = C'' = A$$

For an arbitrary set of objects D the formal concept is computed as follows:

$$D' = A$$

$$A' = D'' = B$$

In total there exists $O(2^{\min(|G|, |M|)})$ itemsets where G and M are the set of objects and attributes, respectively

For every formal concept (A, B) the following holds: $A'' = A$, $B'' = B$

Example. Computing formal concepts

Compute the formal concepts that contain:

1. *attributes*: bread and milk

$\{\text{bread, milk}\}' = \{\mathbf{t1}\},$
 $\{\mathbf{t1}\}' = \{\text{bread, milk}\}'' = \{\mathbf{\text{rice, bread, milk}}\}$

2. *attributes* milk and pasta

$\{\text{milk, pasta}\}' = \{\mathbf{t2, t5}\},$
 $\{\mathbf{t2, t5}\}' = \{\text{milk, pasta}\}'' = \{\mathbf{\text{milk, pasta}}\}$

3. *objects*: t1 and t4

$\{\mathbf{t1, t4}\}' = \{\mathbf{\text{rise}}\}$
 $\{\mathbf{\text{rise}}\}' = \{\mathbf{t1, t4}\}'' = \{\mathbf{\text{t1, t2, t4}}\}$

4. *objects*: t3 and t4

$\{\mathbf{t3, t4}\}' = \{\mathbf{\text{apples}}\}$
 $\{\mathbf{\text{apples}}\}' = \{\mathbf{t3, t4}\}'' = \{\mathbf{\text{t3, t4}}\}$

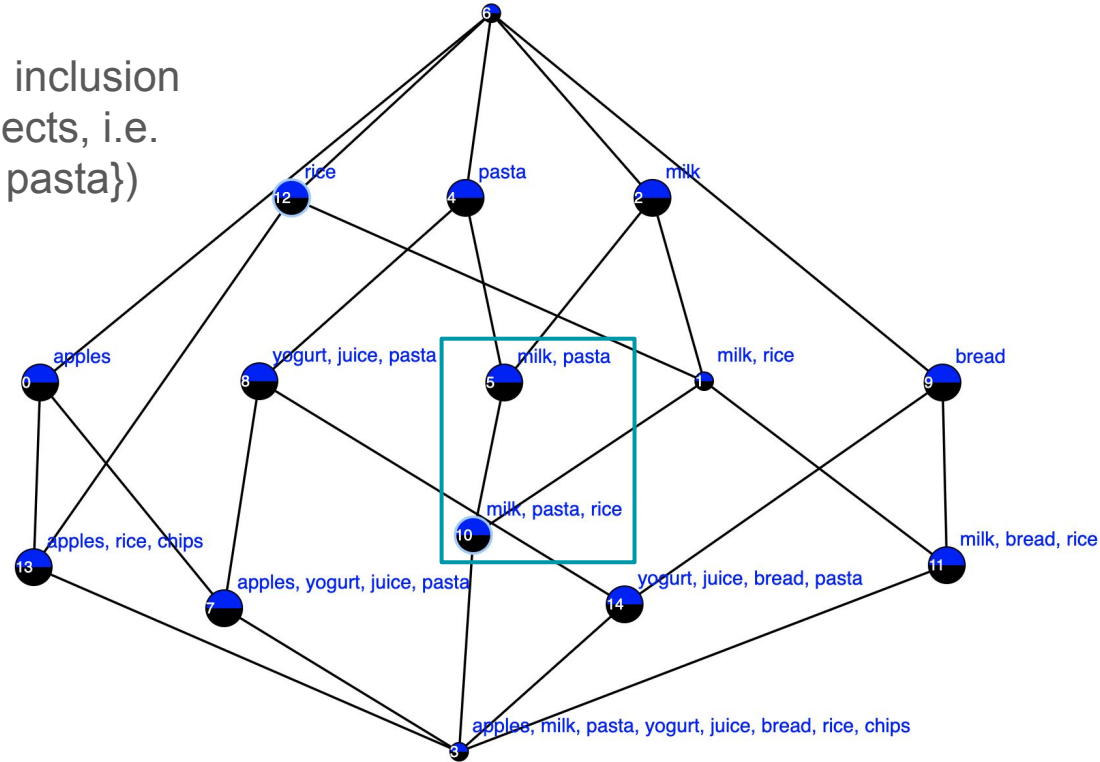
	rice	bread	milk	pasta	yogurt	juice	apples	chips
t1	X	X	X					
t2	X		X	X				
t3				X	X	X	X	
t4	X						X	X
t5			X	X				
t6		X		X	X	X		

Formal concept analysis. Concept lattice

All closed itemsets are ordered by inclusion w.r.t. the set of attributes or a set of objects, i.e.
 $(\{t_2, t_5\}, \{\text{milk}, \text{pasta}\}) > (\{t_2\}, \{\text{rice}, \text{milk}, \text{pasta}\})$

The **concept lattice** contains the whole set of formal concepts (closed itemsets)

If we choose closed itemsets as patterns, the concept lattice is the corresponding pattern search space



Itemset mining

Goal: select a small set of interesting and non-redundant itemsets that cover almost entirely the whole dataset

Standard pipeline (pattern mining):

1. frequent (closed) itemset enumeration
2. application of an interestingness measure to each pattern individually

Drawbacks:

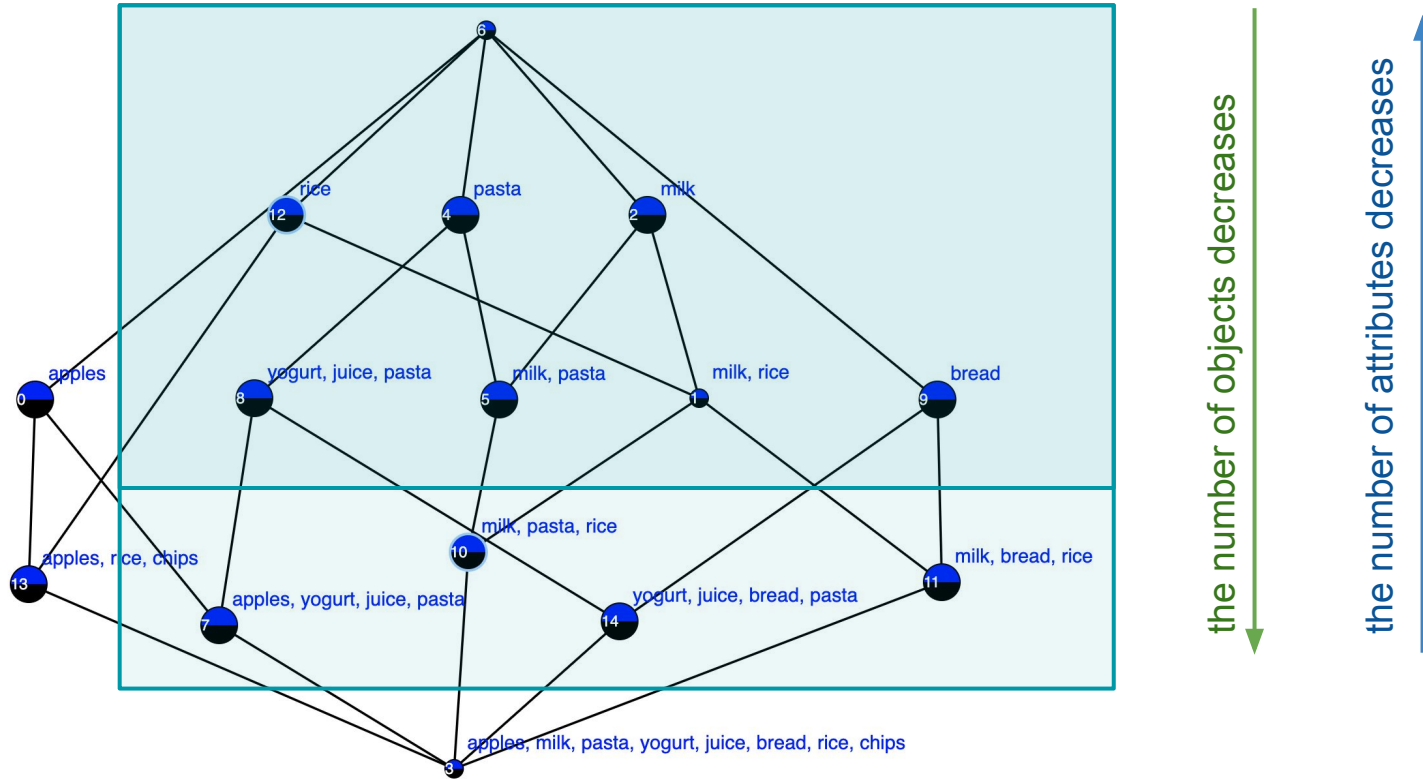
there are no rules of thumb for choosing a good threshold

itemsets are highly redundant

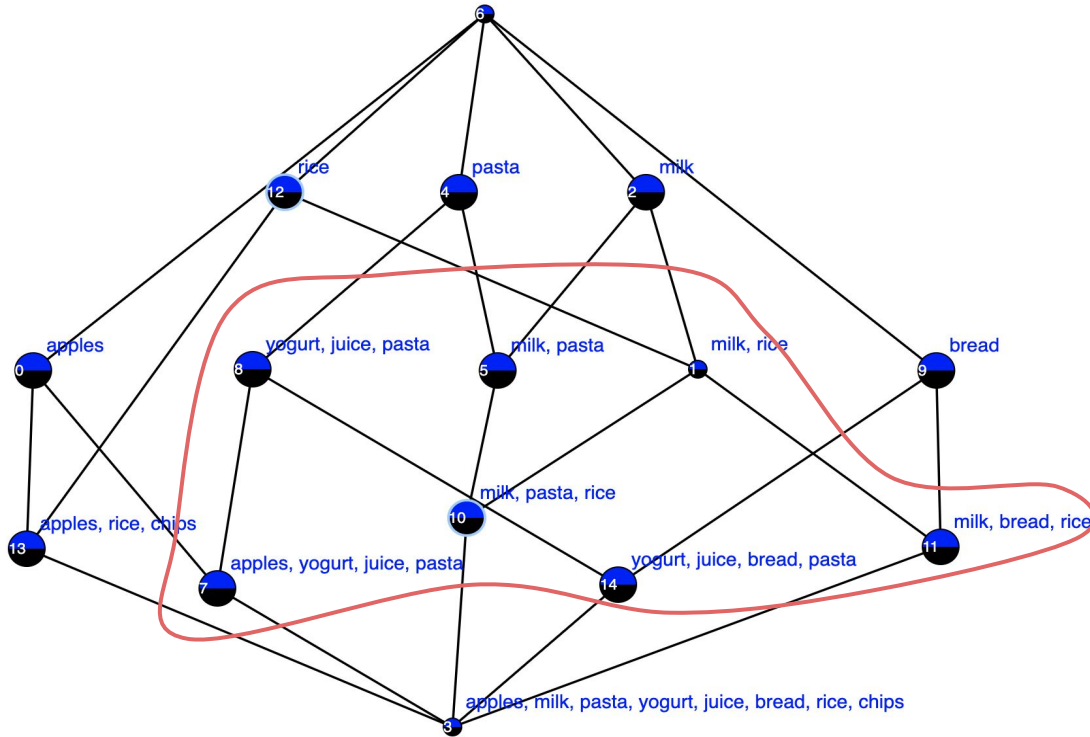
cover is not regulated (usually)

itemsets are not interesting

Enumeration of frequent closed itemsets



Evaluation of almost the same itemsets



Given a set of enumerated itemsets, the problem with application of interestingness measures is the following:

being applied to almost the same itemsets we will see that **almost the same patterns** will obtain **almost the same interestingness scores** =>

very redundant set of interesting itemsets

Example of a set of itemsets

The first set of itemset: {rise, milk}, {bread}, {milk, pasta}, {apples}, {pasta, yogurt, juice}

- covers almost all crosses (cover rate is 17/19)
- is almost non-overlapping (overlapping rate is 18/17)

	rice	bread	milk	pasta	yogurt	juice	apples	chips
t ₁	X	X	X					
t ₂	X		X	X				
t ₃				X	X	X	X	
t ₄	X						X	X
t ₅			X	X				
t ₆		X		X	X	X		

The second set of itemset: {rise}, {bread}, {milk}, {milk, pasta}, {pasta}, {pasta, yogurt, juice}, {rice, apples, chips}

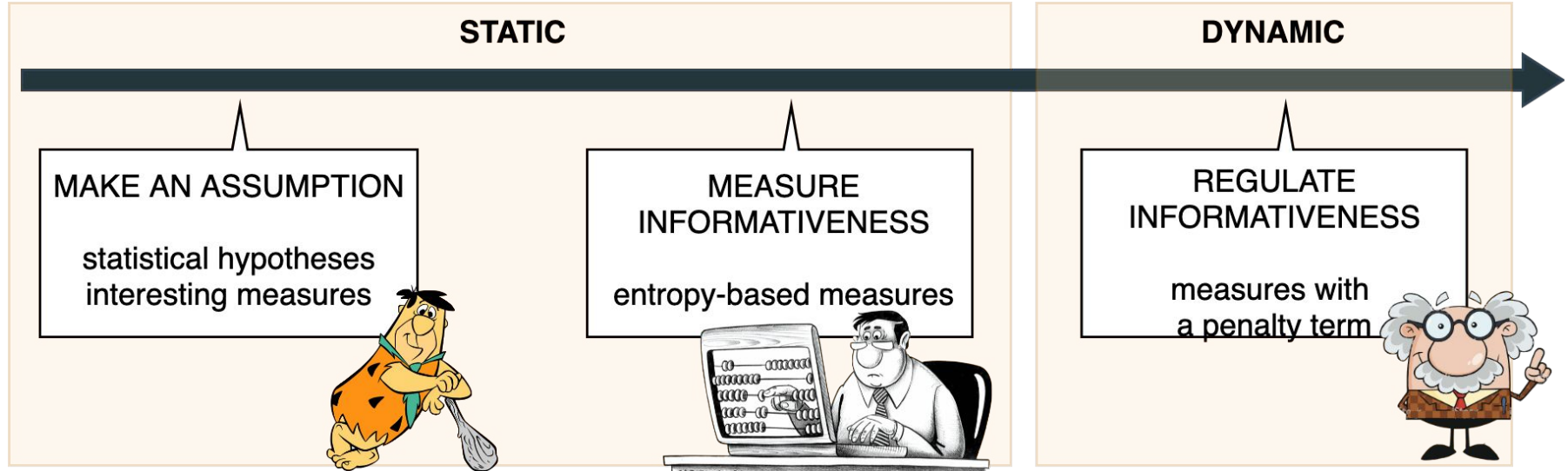
- covers almost all crosses (cover rate is 18/19)
- is almost non-overlapping (overlapping rate is 25/18)

	rice	bread	milk	pasta	yogurt	juice	apples	chips
t ₁	X	X	X					
t ₂	X		X	X				
t ₃				X	X	X	X	
t ₄	X						X	X
t ₅			X	X				
t ₆		X		X	X	X		

Cover rate = #crosses covered at least once / # total crosses

Overlapping rate = # total area of itemsets / # crosses covered at least once

Evolution of itemset mining



Itemset mining

Goal: select a small set of interesting and non-redundant itemsets that cover almost entirely the whole dataset

Standard pipeline (pattern mining):

1. frequent (closed) itemset enumeration
2. application of an interestingness measure to each pattern individually

Drawbacks:

there are no rules of thumb for choosing a good threshold

itemsets are highly redundant

cover is not regulated (usually)

itemsets are not interesting

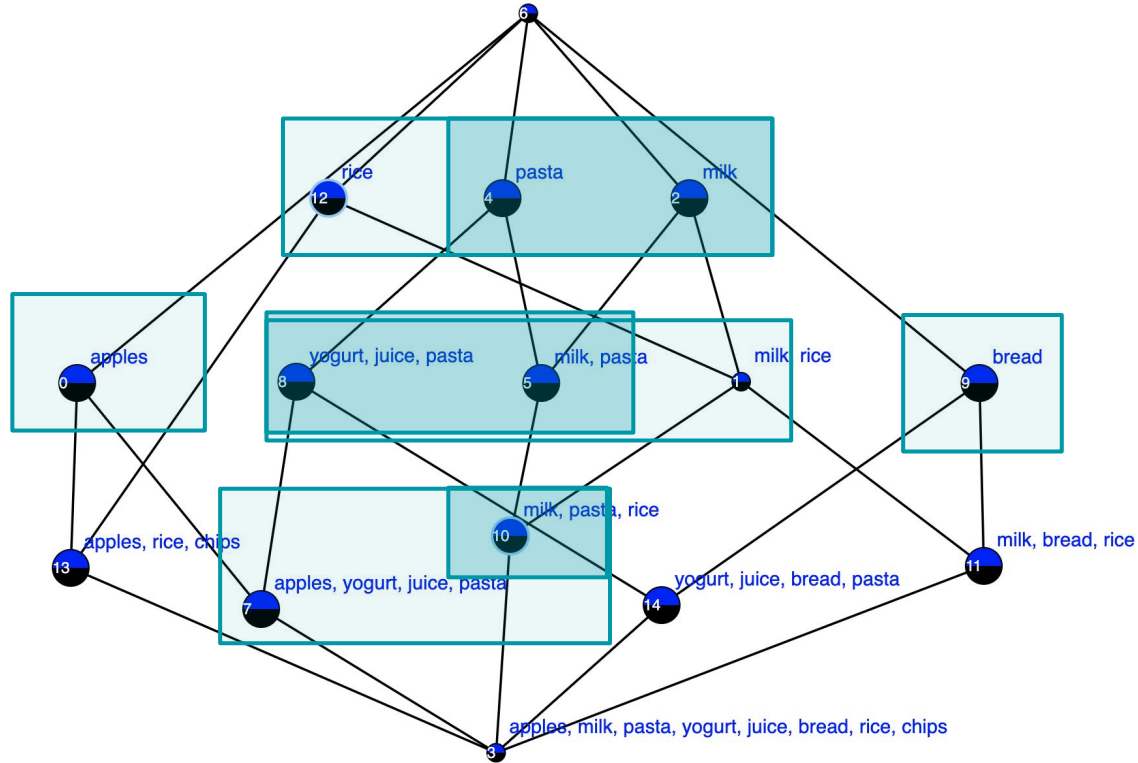
Itemset set mining

Goal: select a small set of interesting and non-redundant itemsets that cover almost entirely the whole dataset

Academic state-of-the-art pipeline (pattern set mining):

- 1.iterative itemset enumeration
- 2.application of a measure to a **set** of itemsets

Formal concept analysis. Concept lattice



Minimum description length principle

MDL-based methods are often referred to as “**learning by compression**”

“Learning” \approx “finding regularities”, thus the more regularities are found, the more the data can be compressed

In MDL the **total description length** is given by

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

where $L(H)$ is the description length of the model (set of patterns) H , in bits, and $L(D|H)$ is the description length of the dataset D encoded with this set of patterns, in bits

Minimum description length principle in pattern mining

Itemsets

- Vreeken, J., Van Leeuwen, M., and Siebes, A. (2011). Krimp: mining itemsets that compress. DMKD, 23(1):169–214.
- Smets, K. and Vreeken, J. (2012). Slim: Directly mining descriptive patterns. ICDM, 236–247. SIAM.

Tiles

- Tatti, N. and Vreeken, J. (2008). Finding good itemsets by packing data. ICDM; 588–597. IEEE.
- Tatti, N. and Vreeken, J. (2012). Discovering Descriptive Tile Trees – By Mining Optimal Geometric Subtiles. ECML PKDD, LNCS 7523, 9–24. Springer

Boolean matrix factorization

- Miettinen, P. and Vreeken, J. (2014). MDL4BMF: Minimum Description Length for Boolean Matrix Factorization. ACM TKDD, 8(4):1–31.

Arbitrary-shaped patterns

- Faas, M. and van Leeuwen, M. (2020) Vouw: Geometric pattern mining using the MDL principle. International Symposium on Intelligent Data Analysis, 158–170. Springer.

Graphs

- Saran, D. and Vreeken, J. (2019). Summarizing dynamic graphs using MDL.

Hyper-rectangles

- Witteveen, J., Duivesteijn, W., Knobbe, A., and Grünwald, P. (2014). Realkrimp – finding hyper intervals that compress with MDL for real-valued data., International Symposium on Intelligent Data Analysis, 368–379. Springer.
- Makhalova, T., Kuznetsov, S. O., & Napoli, A. (2020). Mint: MDL-based approach for Mining INTEResting Numerical Pattern Sets. arXiv preprint arXiv:2011.14843

Discretization

- Nguyen, H.-V., Müller, E., Vreeken, J., and Böhm, K. (2014). Unsupervised interaction-preserving discretization of multivariate data. DMKD, 28(5-6):1366–1397.

MDL in itemset mining. Slim algorithm

In practice, minimization is performed into **two stages**:

- computing pattern candidates
- selection those that minimize the total description length

Main components of the MDL-based approaches:

- encoding scheme
- itemset candidate enumeration
- total description length minimization

Slim. Encoding scheme

apple, bread, duck
apple, bread, duck
apple, bread, carrot
apple, bread, carrot
apple, bread, carrot, eggs, fish
apple, bread, carrot, eggs, fish

encoding with uniform codes				
a	b	d		
a	b	d		
a	b	c		
a	b	c		
a	b	c	e	f
a	b	c	e	f

With uniform encoding, each item is encoded by the code words of uniform length

encoding with ST				
a	b	d		
a	b	d		
a	b	c		
a	b	c		
a	b	c	e	f
a	b	c	e	f

implicit code table	
singletons	code words
a	a
b	b
c	c
d	d
e	e
f	f

We can gain some bits by using the codewords of variable length (frequent items are encoded with shorter codewords). See Shannon codes for details. But we need to store a dictionary that contains the uniform codes and the associated code words of variable length

Slim. Encoding scheme

encoding with ST	standard code table ST		
	singletons	code words	usage
a b d	a	a	6
a b d	b	b	6
a b c	c	c	4
a b c	d	d	2
a b c e f	e	e	2
a b c e f	f	f	2

Applied to itemset mining, the simplest model is one that contains only singletons in its code table.

The goal is to find itemsets that ensure the shortest total description length

abd	
abd	
abc	
abc	
abc	ef
abc	ef

SLIM code table			usage
itemsets	code words		
a b c	abc		4
a b d	abd		2
e f	ef		2

- x variable-length code word associated with the singleton x
- X variable-length code word associated with the itemset X

Slim. Encoding scheme

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

SLIM code table			usage
itemsets	code words		
<div><div>a</div><div>b</div><div>c</div></div>	<div>abc</div>	4	
<div><div>a</div><div>b</div><div>d</div></div>	<div>abd</div>	2	
<div><div>e</div><div>f</div></div>	<div>ef</div>	2	

encoding with the SLIM code table

g_1	<i>abd</i>	
g_2	<i>abd</i>	
g_3	<i>abc</i>	
g_4	<i>abc</i>	
g_5	<i>abc</i>	<i>ef</i>
g_6	<i>abc</i>	<i>ef</i>

$$L(CT) = \underbrace{2 \cdot (l(a) + l(b)) + l(c) + l(d) + l(e) + l(f)}_{\substack{L(abd|ST) + L(abc|ST) + L(ef|ST) \\ \text{the left-hand column length}}} + \underbrace{L(\text{code}(abc)) + L(\text{code}(abd)) + L(\text{code}(ef))}_{\text{the right-hand column length}}.$$

$$L(D|CT) = 2 \cdot (L(\text{code}(abd)) + L(\text{code}(ef))) + 4 \cdot L(\text{code}(abc)),$$

Slim. Candidate discovery by pairwise merging

Step 1	Candidates	Step 2	Candidates	Step 3	Candidates	Step 4	Candidates	Step 5	Cand.
	<i>b c d e f</i>		<i>c d e f</i>		<i>ef c d</i>		<i>ef d</i>		<i>ef</i>
<i>a</i>	<i>ab</i> <i>ac ad ae af</i>	<i>ab</i> <i>abc abd abe abf</i>	<i>ab</i> <i>abef</i> <i>abc</i> <i>abd</i>	<i>abc</i> <i>abcdef</i> <i>abcd</i>	<i>abc</i> <i>abcdef</i>				
<i>b</i>	<i>bc bd be bf</i>	<i>c</i> <i>cd ce cf</i>	<i>ef</i> <i>cef cdef</i>	<i>ab</i> <i>abef</i> <i>abd</i>	<i>abd</i>				
<i>c</i>	<i>cd ce cf</i>	<i>d</i> <i>de df</i>	<i>c</i> <i>cd</i>	<i>ef</i> <i>def</i>					
<i>d</i>	<i>de df</i>	<i>e</i> <i>ef</i>	<i>d</i>						
<i>e</i>	<i>ef</i>	<i>f</i>							
<i>f</i>									

At each step, the candidates to the code table are computed by merging the members of a current code table.

In bold, the patterns appearing for the first time in the code table or in the candidate set are highlighted.

patterns from the code table at the current step,
 a set of candidates
 selected pattern

Association rule mining vs classification rule mining

transactions
apple, bread, duck
apple, bread, duck
apple, bread, carrot
apple, bread, carrot
apple, bread, carrot, eggs, fish
apple, bread, carrot, eggs, fish

$d \rightarrow ab$

$e \rightarrow f$

$cf \rightarrow abe$

$ab \rightarrow c$

$ab \rightarrow d$

$d \rightarrow 0$

$c \rightarrow 1$

$cf \rightarrow 1$

Association rule mining

If **antecedent** then **consequent** : **A** \rightarrow **C**

- $\text{support}(A \rightarrow C) = \text{support}(A \cup C) = P(AC) \in [0,1]$
- $\text{confidence}(A \rightarrow C) = \text{support}(A \rightarrow C) / \text{support}(A) = P(C|A) \in [0,1]$
- $\text{lift}(A \rightarrow C) = \text{confidence}(A \rightarrow C) / \text{support}(C) = P(AC) / (P(A)P(C)) \in [0,\text{inf})$
- $\text{leverage}(A \rightarrow C) = \text{support}(A \rightarrow C) - \text{support}(A)\text{support}(C) = P(AC) - P(A)P(C) \in [-1,1]$

Example. Computing the best rules

transactions
apple, bread, duck
apple, bread, duck
apple, bread, carrot
apple, bread, carrot
apple, bread, carrot, eggs, fish
apple, bread, carrot, eggs, fish

measures	$c \rightarrow ef$
$\text{support}(c \rightarrow ef) = P(cef)$	$2 / 6$
$\text{confidence}(c \rightarrow ef) = P(ef \mid c)$	$2 / 4$
$\text{lift}(c \rightarrow ef) = P(cef) / (P(c)P(ef))$	$2 * 6 / (2 * 4) = 12 / 8$
$\text{levarage}(c \rightarrow ef) = P(cef) - P(c)P(ef)$	$2 / 6 - 4 / 6 * 2 / 6 =$ $= 1 / 9$

Reminder:

$$\text{support}(A \rightarrow C) = P(AC) \in [0,1]$$

$$\text{confidence}(A \rightarrow C) = P(AC) / P(A) = P(C|A)$$

$$\text{lift}(A \rightarrow C) = P(AC) / (P(A)P(C))$$

$$\text{levarage}(A \rightarrow C) = P(AC) - P(A)P(C)$$

Causal relations in association rules

Associations may not imply causal relationships!

A measure of the risk of experiencing the outcome under study when the antecedent factor is present is given by

$$\Omega_A = \frac{P(C|A)}{P(\bar{C}|A)}$$

Ω_A is the odds A that B will occur when A is present. The odds $\Omega_{\neg A}$ are defined similarly

$$\Omega_{\bar{A}} = \frac{P(C|\bar{A})}{P(\bar{C}|\bar{A})}$$

The odds ratio is given by

$$\omega_D = \frac{\Omega_A}{\Omega_{\bar{A}}} = \frac{\text{supp}(AC)\text{supp}(\bar{A}\bar{C})}{\text{supp}(\bar{A}C)\text{supp}(A\bar{C})}$$

Interpretation: $\omega_D = 1$ - exposure A does not affect odds of outcome C
 $\omega_D > 1$ - exposure A associated with higher odds of outcome C
 $\omega_D < 1$ - exposure A associated with lower odds of outcome C

Example. Computing the best rules

transactions	$A \rightarrow C$	$c \rightarrow ef$
apple, bread, duck	$P(C A) = P(AC)/P(A)$ $P(\neg C A)$ $P(C \neg A)$ $P(\neg C \neg A)$ $\Omega_A = P(C A) / P(\neg C A)$ $\Omega_{\neg A} = P(C \neg A) / P(\neg C \neg A)$ $\omega = \Omega_A / \Omega_{\neg A}$	
apple, bread, duck		
apple, bread, carrot		
apple, bread, carrot		
apple, bread, carrot, eggs, fish		
apple, bread, carrot, eggs, fish		

Example. Computing the best rules

transactions	$A \rightarrow C$	$c \rightarrow ef$
apple, bread, duck	$P(C A) = P(AC)/P(A)$	2 / 4
apple, bread, duck	$P(\neg C A)$	2 / 4
apple, bread, carrot	$P(C \neg A)$	0
apple, bread, carrot	$P(\neg C \neg A)$	2 / 2 = 1
apple, bread, carrot	$\Omega_A = P(C A) / P(\neg C A)$	1
apple, bread, carrot, eggs, fish	$\Omega_{\neg A} = P(C \neg A) / P(\neg C \neg A)$	0
apple, bread, carrot, eggs, fish	$\omega = \Omega_A / \Omega_{\neg A}$	0

Simpson paradox

	Salary = <i>low</i>	Salary = <i>high</i>
Gender = <i>m</i>	185	120
Gender = <i>f</i>	65	60

$$\omega_D(m \rightarrow low) = \omega_D(w \rightarrow high) = \frac{185 \cdot 60}{65 \cdot 120} = 1.42$$

The odds ratio indicates a positive association between “Gender = m” and “Salary = low”

	Salary = <i>low</i>	Salary = <i>high</i>
Gender = <i>m</i> & College = <i>y</i>	5	20
Gender = <i>f</i> & College = <i>y</i>	15	40

$$\begin{aligned} \omega_{D_{college=y}}(m \rightarrow low) &= \omega_D(w \rightarrow high) \\ &= \frac{5 \cdot 40}{15 \cdot 20} = 0.66 \end{aligned}$$

	Salary = <i>low</i>	Salary = <i>high</i>
Gender = <i>m</i> & College = <i>n</i>	180	100
Gender = <i>f</i> & College = <i>n</i>	50	20

$$\begin{aligned} \omega_{D_{college=n}}(m \rightarrow low) &= \omega_D(w \rightarrow high) \\ &= \frac{180 \cdot 20}{50 \cdot 100} = 0.72 \end{aligned}$$

Studying the causal relations

Studying the causal relations like $A \rightarrow C$, it is important to consider two subsets (groups), with A and without A, such that the instances within these two groups are the most **similar** by **relevant** to C characteristics

Difficulties:

- *Similarity*: What is similarity? Is there any threshold on similarity?
- *Relevance*: Which attributes are relevant for given C?

Example: “man” \rightarrow “low salary”

Similarity

Can we consider the age of 25 and 30 years to be similar?

What about 25 and 35? 35 and 45?

Relevance

Possible relevant characteristic: age, education, domain

Irrelevant characteristics: hair and eyes color, food preferences

Numerical pattern mining

Goal: to discover a small set of diverse and interesting patterns that cover almost entirely the whole dataset

Pattern is a **subset of the input feature space** that contain similar objects

The chosen type of patterns affect the interpretability of the results. Patterns of “simple” shape may have a limited descriptiveness

Issues:

- Computationally **expensive search for splitting points** in each interval range (compare with a greedy dichotomous splitting of the decision trees)
- **Pre-discretization** is not always a good solution (in a few slides)
- **Potentially exponential number of patterns**, the approaches for a greedy search for good rules are underdeveloped

Pattern shape

Selecting a particular shape of patterns (pattern type) we may

- restrict the pattern search space
- improve the interpretability of the results

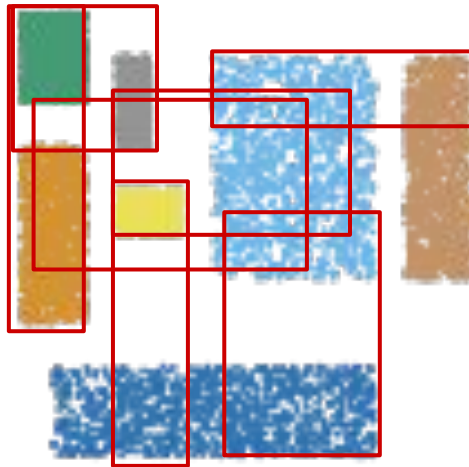
The best pattern shapes w.r.t. the aforementioned goals are

- spheres (are set by their center and radius)
- hyper-rectangles (are set by a tuple of intervals)

Pattern search space

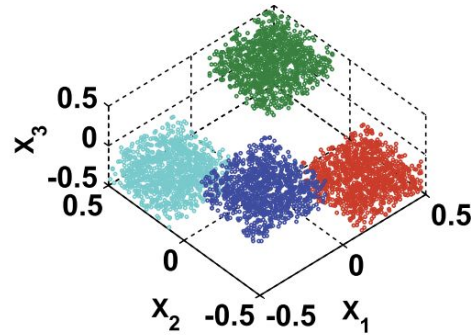
For hyper-rectangles the pattern search space is enormous

How to reduce the pattern search space?

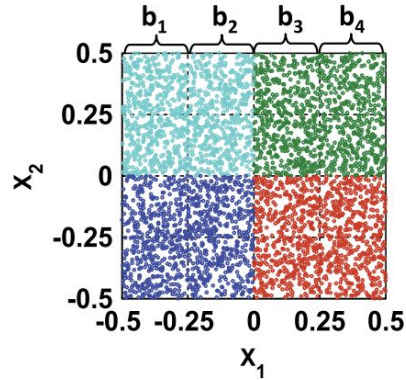


Let's **discretize** the search space! But how to do it **correctly**?

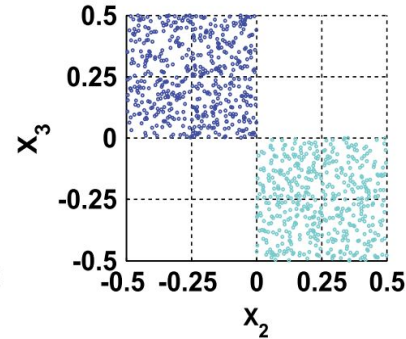
Attribute-wise discretization



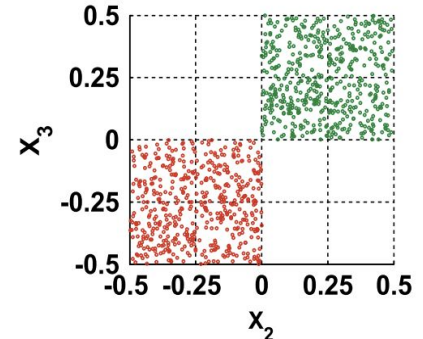
(a) (X_1, X_2, X_3)



(b) (X_1, X_2)



(c) $p(X_2, X_3 \mid b_1)$ and
 $p(X_2, X_3 \mid b_2)$

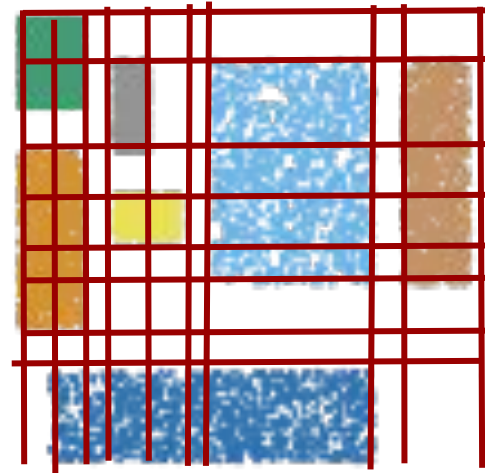
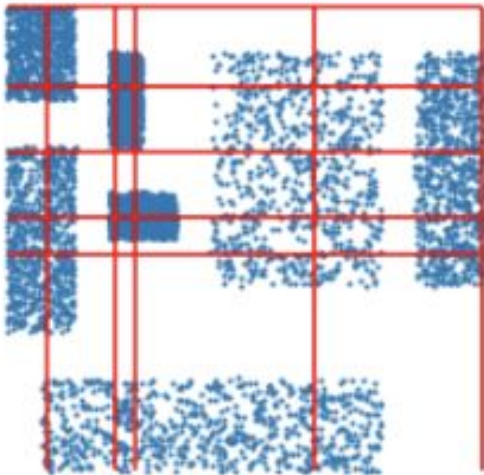


(d) $p(X_2, X_3 \mid b_3)$ and
 $p(X_2, X_3 \mid b_4)$

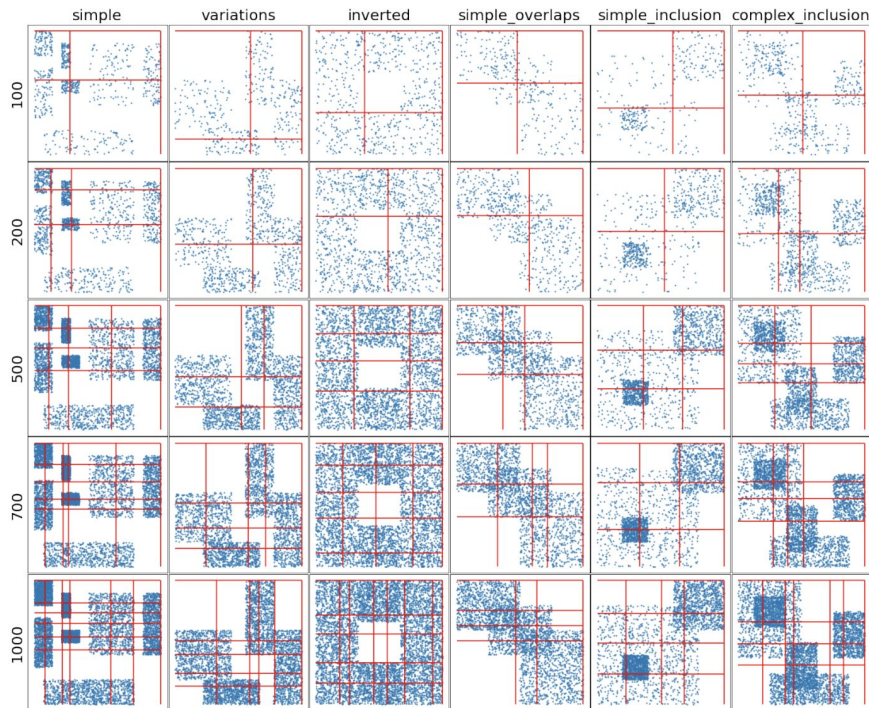
Attributes should be discretized w.r.t. other attributes to preserve interaction!

Uniform vs pattern-specific discretization

An “optimal” discretization is one that **preserves interactions** between features and put enough boundaries (not too much and not too little) in the feature space



Coarse vs pattern-specific discretization



Finding globally optimal interaction-preserving splitting does not allow to get patterns with precise boundaries

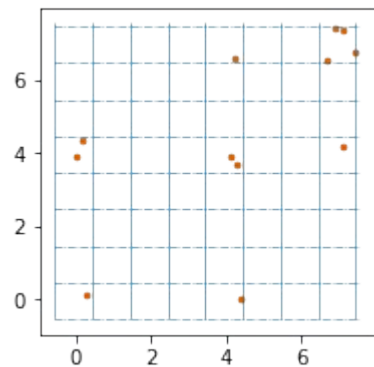
Makhalova, Tatiana, Sergei O. Kuznetsov, and Amedeo Napoli. "Mint: MDL-based approach for Mining INteresting Numerical Pattern Sets." arXiv preprint arXiv:2011.14843 (2020).

Boundaries are computed by IPD, source code: <https://people.mmci.uni-saarland.de/~jilles/prj/ipd/>

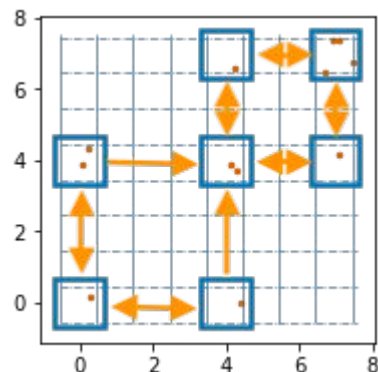
MDL for numerical pattern mining

- Use a fine-grained discretization
- Define an encoding scheme
- Gradually mine candidates and select those that ensure the best compression

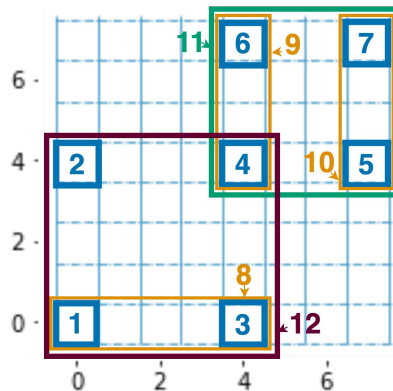
Example:



fine-grained discretization



a candidate is a smallest rectangle that contain two closest patterns



gradually select those candidates that ensure the best compression

Encoding of MINT

$$L(D, \mathcal{H}) = L(\mathcal{H}) + L(D|\mathcal{H})$$

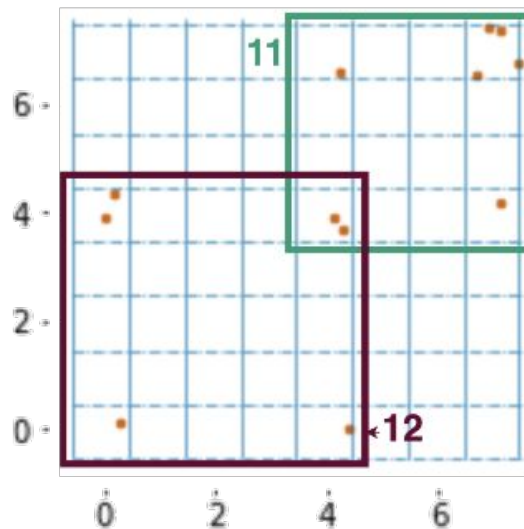
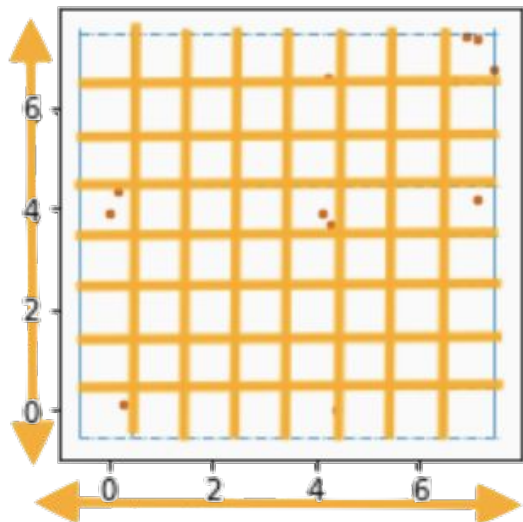
$$L(\mathcal{H}) = \underbrace{L_{\mathbb{N}}(|M|) + \sum_{i=1}^{|M|} L_{\mathbb{N}}(|\mathcal{B}_i|)}_{\text{length of the grid}} + \underbrace{L_{\mathbb{N}}(|\mathcal{H}|) + |\mathcal{H}| \left(\sum_{i=1}^{|M|} \log (|\mathcal{B}_i|(|\mathcal{B}_i| + 1)/2) \right)}_{\text{length of the pattern set}}$$

$$L(D|\mathcal{H}) = \underbrace{L_{\mathbb{N}}(|G|)}_{\text{cost of encoding the number of instances}} + \underbrace{L(D(\mathcal{H})|\mathcal{H})}_{\text{cost of encoding } D \text{ with } \mathcal{H}} + \underbrace{L(D \ominus D(\mathcal{H})|\mathcal{H})}_{\text{reconstruction cost}}$$

Length of the model (pattern set)

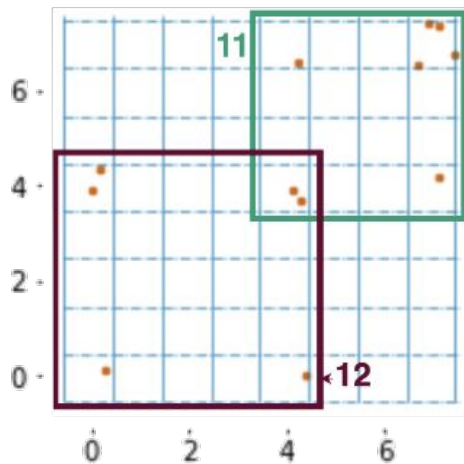
$$L(\mathcal{H}) = L_{\mathbb{N}}(|M|) + \sum_{i=1}^{|M|} L_{\mathbb{N}}(|\mathcal{B}_i|) + L_{\mathbb{N}}(|\mathcal{H}|) + |\mathcal{H}| \left(\sum_{i=1}^{|M|} \log(|\mathcal{B}_i|(|\mathcal{B}_i| + 1)/2) \right)$$

$$L(\mathcal{H}) = \underbrace{L_{\mathbb{N}}(2) + L_{\mathbb{N}}(8) + L_{\mathbb{N}}(8)}_{\text{length of the grid}} + \underbrace{L_{\mathbb{N}}(2) + 2(\log(36) + \log(36))}_{\text{length of the pattern set}}$$



Length of data

$$L(D|\mathcal{H}) = \underbrace{L_{\mathbb{N}}(12)}_{\text{cost of encoding the number of instances}} + \underbrace{L(D(\mathcal{H})|\mathcal{H})}_{\text{cost of encoding } D \text{ with } \mathcal{H}} + \underbrace{L(D \ominus D(\mathcal{H})|\mathcal{H})}_{\text{reconstruction cost}}$$



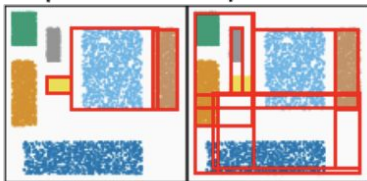
$$\text{usg}(h_{11}) = 8$$

$$\text{usg}(h_{12}) = 4$$

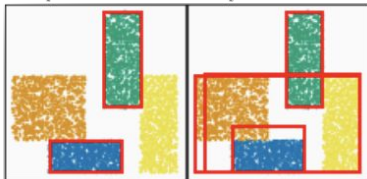
Comparison of pattern miners

REALKRIMP

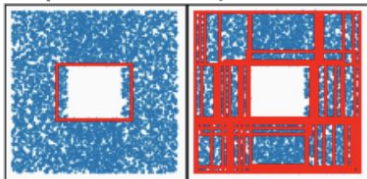
#pat. = 3 #pat. = 9



#pat. = 2 #pat. = 5

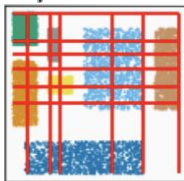


#pat. = 1 #pat. = 47

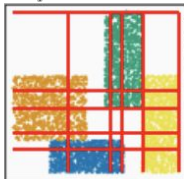


IPD

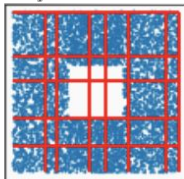
#pat. = 24



#pat. = 24

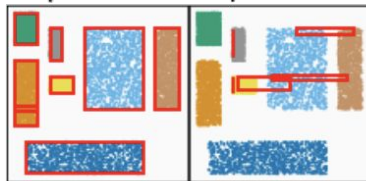


#pat. = 24

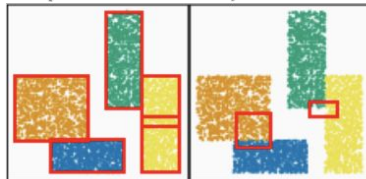


MINT

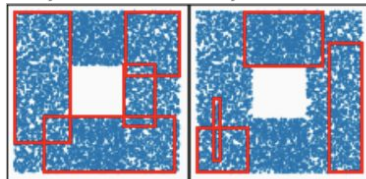
#pat. = 8 #pat. = 5



#pat. = 5 #pat. = 2



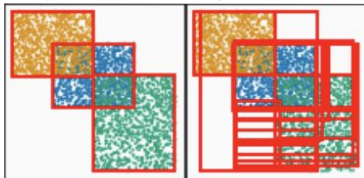
#pat. = 4 #pat. = 4



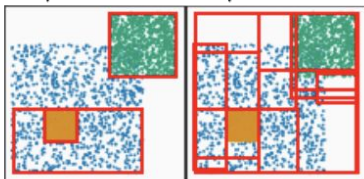
Comparison of pattern miners

REALKRIMP

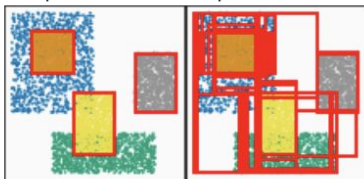
#pat. = 3 #pat. = 18



#pat. = 3 #pat. = 9

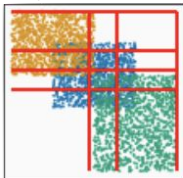


#pat. = 3 #pat. = 38

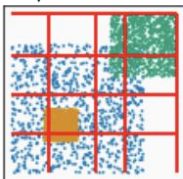


IPD

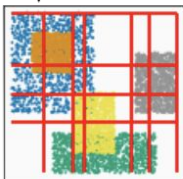
#pat. = 24



#pat. = 24

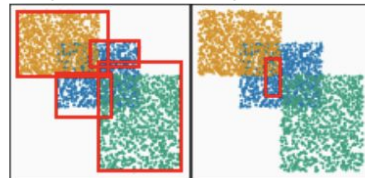


#pat. = 24

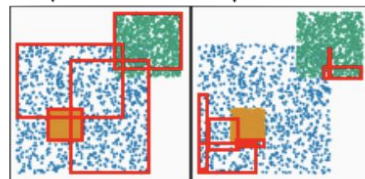


MINT

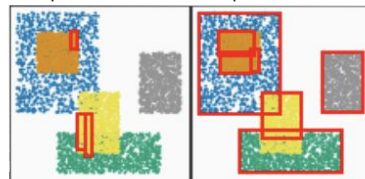
#pat. = 4 #pat. = 1



#pat. = 4 #pat. = 6



#pat. = 3 #pat. = 7



Mint: pros and cons

Pros:

- returns diverse and non-redundant patterns
- patterns have precise boundaries
- pattern descriptions are simple -> the results are interpretable

Cons:

- is based under the assumption on uniform distribution
- may be sensitive to noise
- requires a lot of memories
- does not contain “excluded” patterns in its language

Numerical pattern mining vs clustering

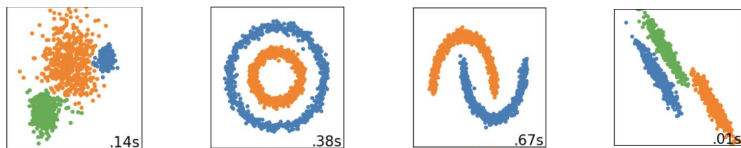
In pattern mining (PM) the **description** comes first, while in clustering the primacy is given to the **object similarity**

PM entails some requirements for the ease of **interpretation**, i.e., the resulting patterns should describe a region in the “attribute space” that is easy to interpret. In clustering, the focus is put on **groups** of objects or instances

The clusters can be constrained to have certain shapes, e.g., spheres in K-means or DBSCAN, but still the similarity of objects remains the most important characteristic of clusters

Numerical pattern mining vs clustering

Clustering is better when we need to identify groups of a complex shape but do not need to represent them in the attribute space



In the contrast to numerical pattern mining, clustering has the following **drawbacks**:

- selecting a number of clusters
- selecting a distance metric
- dealing with the curse of dimensionality
- uninterpretable description of patterns

Wrapping up

- The state-of-the-art pattern mining approaches
 - discover pattern search space iteratively
 - evaluate pattern set instead of each pattern individually, e.g., using MDL
- Association rules (antecedent → consequent)
 - do not say anything about causal relations
 - can be considered as classification rules (if consequent is a class label)
- Pre-discretization is not a good idea (should be avoided if possible)