# The Data Problem in Data Mining

#### Albrecht Zimmermann Université de Caen Normandie

#### 14th IDA, St. Étienne, 24.10.2015

UNIVERSITÉ CAEN NORMANDIE

### The Data Problem in Unsupervised Pattern Mining

#### Albrecht Zimmermann Université de Caen Normandie

#### 14th IDA, St. Étienne, 24.10.2015

UNIVERSITÉ CAEN NORMANDIE

### Take home message

#### We don't understand

- Algorithm run times
- Parameter settings
- How to interpret mined patterns

#### We can begin to fix this

- Researching data generation
- Generating data
- Exploring mining behavior

→Because we lack data !

Add new knowledge/tools

## Local pattern mining









### Road map

### 1) Run time behavior

- 2) Parameter setting/output behavior
- 3) Relation patterns data
- 4) Data science

### Road map

## 1) Run time behavior

- Problem setting
- Existing evaluation
- Data problem
- Attempts at understanding







# Run times – how much data ?



# Run times – how much data ?



### What's the issue ?

Not enough data !

Zheng data sets added

-Nowadays FIMI (Goethals & Zaki) repository 12 data sets

UCI : 333 UCR : 70+

Zheng et al. killed QUEST, no replacement

# Understanding run times

Pattern type	Itemsets	Graphs	Episodes
Characterized by	Distribution pattern lengths	<ul> <li>Subtasks</li> <li>Run time per fragment</li> <li>Edge densities/ Branching factor</li> <li>Java VM/Cache size</li> </ul>	<ul> <li>Length</li> <li>Non-event probability</li> <li>Alphabet size</li> </ul>
	Zaki et al. 2002 Ramesh et al. 2003 Gouda et al. 2005 Flouvat et al. 2010	Wörlein et al. 2005 Nijssen et al. 2006	Zimmermann 2014

# Understanding run times

Pattern type	Itemsets	Graphs	Episodes
Characterized by	Distribution pattern lengths	<ul> <li>Subtasks</li> <li>Run time per fragment</li> <li>Edge densities/ Branching factor</li> <li>Java VM/Cache size</li> </ul>	<ul> <li>Length</li> <li>Non-event probability</li> <li>Alphabet size</li> </ul>
Zaki !	Zaki et al. 2002 Ramesh et al. 2003 Gouda et al. 2005 Flouvat et al. 2010	Wörlein et al. 2005 Nijssen et al. 2006	Zimmermann 2014

### Predicting run times

Based on	Bernoulli/ Markovian model	Partial mining results	Sampling
	Lhote et al. 2005	Palmerini et al. 2004 Geerts et al. 2005	Boley et al. 2008/2010
		Work-around : how many frequent sets?	2





Top-k Pattern set mining Statistical pattern mining MaxEnt



Istical pattern min MaxEnt



- Same data (not much)
- W/problematic properties
- No patterns known !
- Can't compare

# Understanding output sizes/composition

Pattern type	Itemsets	Strings		
Characterized by	Samples + Background models	Samples + Background models		
	Boley et al. 2008/2010 Van Leeuwen et al. 2014	Besson et al. 2008		
	Work-a where spectrum from ra	around : e's the n different andom ?		

Dairy	Produce	Produce	Flowers		Customer Service	
Dairy	Coffee	Frozen Foods Canned Foods Tea	Cercal	Promo	Registers	
Dany	Boor	Bakery Snack Foods Wine	Beverages	Water		Coffee Grinder









## Understanding patterns – evaluation

Ad-hoc recovery of embedded patterns

Flanking run time experiments
Consulting domain experts

- Arguably biased

Humans see patterns everywhere !



### No data !

- ...w/ known ground truth
- Prediction has labels
- Minimum : known patterns
- Better : known processes

## Understanding patterndata relation

Ground Truth	Embedded patterns	Generative models		
	Zimmermann 2013, 2014	Mampaey et al. 2013 Webb et al. 2014		
	Comparing found to embedded	Comparing found to substructures of model		

# Proposal : generating data

### QUEST led the way

- -Allows for different characteristics
- -Known patterns
- -Controllable properties
- Others took it up
- -Non-systematically

We wouldn't be the first Engineering

- -Downs et al. 1993
- -CORSIKA Physics
- -Analytical sociology
- -Pei et al. 2006
- -SAT solving
- -Pascal this morning

QUEST!

Wrong transaction length distribution Replacement ? Cooper 2009 ?



log(support)

**QUEST! Replacement**? Wrong transaction Cooper 2009 ? length distribution 100 10000 QUEST 15k items Retail 1000 og([{i | support(i) = x}]) og(|{i | support(i) = x}) 10 100 10 1 10 100 1000 10000 1 10 100 1000 1 10000

100000

log(support)





Wrong transaction length distribution

### **Replacement**? Cooper 2009 ?

Adä et al. The new iris data: modular data generators. KDD 2010

#### Inverse itemset No patterns mining **Depend on**

QUEST!

existing data

### Chakrabarti et al.

2006

20 graph generators None gets it right No patterns

## Challenge : recoverable patterns/realistic processes

- Sampling a degree distribution is **not** realistic !
- -And one cannot compare a subgraph to it
- Patterns don't have to be the same as data
- -Bayes' Nets to build itemsets
- -Sequences to build graphs
- -Agents !

## Challenge : describing/ comparing data

When are data sets similar ?

The MetaLearners know about this :)

- -Column/row marginals not enough
- -Density not enough
- -Krimp code tables sunk by relabeling
- -Comparisons for structured data ?
- Should similar data lead to similar results ?

# Challenge : making the right assumptions

Need to talk to specialists in the field

- -Who might lack the full picture
- -Who might be wrong
- Patterns might (seem to) make no sense

We'll have to go back time and again

CORSIKA is work in progress...











### What do you think ?

*The Data Problem in Data Mining.* SIGKDD Explorations 16 (2)

## •This slide intentionally left blank