# Dense Neighborhood Pattern Sampling in Numerical Data

**Arnaud Giacometti and Arnaud Soulet** 

**University of Tours, France** 





Neighborhood pattern sampling problem

**Three-step random procedure** 

**Experimental study on UCI benchmarks** 





#### Finding interesting patterns in all subspaces

(pattern = an anomalously high local density of data points [Hand, 2002])







X

y



**Discretized patterns** 



Subspace patterns



**Interval patterns** 

loss of information





#### loss of information + combinatorial explosion





loss of information + combinatorial explosion + curse of dimensionality



### Neighborhood patterns (for contrast mining) [Konijn et al., PAKDD13]





### Limit of neighborhood patterns [Konijn et al., PAKDD13]





# Challenge of neighborhood patterns = infinite search space





# Neighborhood pattern sampling problem

Return a random neighborhood pattern with a probability proportional to its number of neighbors (given a radius)





# Neighborhood pattern sampling problem

Return a random neighborhood pattern with a probability proportional to its number of neighbors density (given a radius)



Normalizing by the volume for dealing with different subspaces



# Neighborhood pattern sampling problem

Return a random neighborhood pattern with a probability proportional to its number of neighbors density (given a radius)



Normalizing by the volume for dealing with different **p-norms** 



# Methods for pattern sampling

#### **Stochastic procedure**

- □ Output Space Sampling for Graph Patterns [Hasan et Zaki, VLDB09]
- Fast Query Execution for Retrieval Models Based on Path-Constrained Random Walks [Lao et Cohen, KDD10]

#### **Two-step random procedure**

- Direct local pattern sampling by efficient two-step random procedures [Boley et al., KDD11]
- Instant Exceptional Model Mining Using Weighted Controlled Pattern Sampling [Moens et Boley, IDA14]

#### SAT procedure

□ Flexible constrained sampling with guarantees for pattern mining [Dzyuba et al., DAMI17]

### Two-step random procedure [Boley et al., KDD11]



Frequent itemset sampling: return a random itemset with a probability proportional to its frequency

(AB has twice more chance to be drawn than AC.)



### Two-step random procedure [Boley et al., KDD11]



**O** Calculate the number of itemsets per transaction



### Two-step random procedure [Boley et al., KDD11]



**①** Draw a transaction  $t \in D$  with a probability proportional to the number of itemsets contained in t





**①** Draw a transaction  $t \in D$  with a probability proportional to the number of itemsets contained in t

**2** Draw uniformly an itemset from *t* 

### Three-step random procedure (1)





### Three-step random procedure (2)



### Three-step random procedure (3)











01/06/2018

### Contributions of three-step random procedure

Exact sampling method of neighborhood patterns with a probability proportional to its density (given a radius and a p-norm)

Generalization of the two-step random procedure (mixed data)

#### **Third step is not costly:**

- I-norm: O(log(#data points) + #dimensions x log(#dimensions))
- 2-norm and ∞-norm: O(log(#data points) + #dimensions)
- In practice, average time per pattern: 70 ms

# Experimental study on UCI benchmarks

#### Protocol

 $\Box$  Z-score for numerical data  $\Box$  Radius = 1

#### Questions

- 1. What is the proportion of patterns that we do not observe in a randomized dataset? (= plausibility)
- 2. What is the proportion of distinct patterns? (= diversity)
- 3. What is the accuracy of a sample-based associative classifier?



### 1. Plausibility of sampled patterns

**QCM-BioChem - Paris** 



**Plausibility =** proportion of patterns that we do not observe in a randomized dataset





01/06/2018

## 2. Diversity of sampled patterns



**Diversity** = proportion of distinct patterns

- Diversity enhancement thanks to the third step
- Diversity as good as interval pattern mining

## 3. Sampling-based associative classification



CBA-like method with a sample of 10k patterns

- Accuracy as good as classic methods with a complete extraction
- Sampling-based classifier built in seconds

### Conclusion

#### Description of the second s

- Instant discovery of patterns despite an infinite search space
- High quality patterns
- Lossless pattern mining

#### □ Next steps:

- Interactive pattern mining
- Instant subspace clustering







# Thank you!

