Efficiently Sampling Interval Patterns from Numerical Databases

Djawad Bekkoucha¹, Lamine Diop², Abdelkader Ouali¹, Bruno Crémilleux¹ and Patrice Boizumault¹

¹ University of Caen Normandy, ENSICAEN, CNRS, Normandy Univ GREYC UMR6072, F-14000 Caen, France ² EPITA Research Laboratory (LRE), Le Kremlin-Bicêtre, Paris, FR-94276, France

djawad.bekkoucha@unicaen.fr

GREYC **Electronics and Computer**

Science Laboratory

1. Approaches for Pattern Discovery



No sampling approach for

interval patterns

2. How to mine patterns from numerical data?

- Numerical databases are widely used in a large number of fields
- Use the interval pattern language to represent informations from numerical data:

Let $\mathcal{V} = \langle [a_i, b_i]_{i \in \{1, \dots, |\mathcal{M}|\}} \rangle$, $a_i, b_i \in \mathcal{N}_i \land a_i \leq b_i$ an interval pattern

	Height	Weight	Age
	m_1	m_2	m_3
g_1	155	74	80
g_2	176	99	74
g_3	167	73	28
g_4	153	76	52
g_5	190	99	76

Example

- ► $cover(\langle [153,155], [73,76], [52,80] \rangle) = \{ g_1, g_4 \}$
- ► $freq(\langle [153,155], [73,76], [52,80] \rangle) = |\{ g_1, g_4 \}| = 2$
- $\blacktriangleright desc(\{ g_1, g_4 \}) = \langle [153, 155], [73, 76], [52, 80] \rangle$

3. State of the Art

4. Frequency-Based Interval Pattern Sampling

Stochastic methods:

- Graphs: [El Hasan et al.VLDB 2009], [Bendimerad et al., ICDM 2016]
- Formal concepts: [Boley et al., DSM 2010]
- **Tiles:** [Bendimerad et al., IDA 2020]

- Goal: compute the sum of frequencies of all interval patterns in search space without enumeration
- Equivalent to counting the number of patterns covering each object $g \in \mathcal{G}$:
- $NIP(g) = \prod |\mathcal{U}(v_{g,m})| + |\mathcal{A}(v_{g,m})| + |\mathcal{U}(v_{g,m})| \cdot |\mathcal{A}(v_{g,m})| + 1$ $\mathcal{U}(v_{g,m}) = \{ v \in \mathcal{N}_m \mid v < v_g \}$ such that $\mathcal{A}(v_{q,m}) = \{ v \in \mathcal{N}_m \mid v > v_q \}$

Pre-processing

153

190

73 28

52

76

99

Example: Let us consider attribute m1 and object g1. The value 155 can be:

- ► A lower bound : [155, 167], [155, 176], [155, 190]
- An upper bound : [153, 155]
- ▶ Strictly included : [153, 167], [153, 176], [153, 190]
- ▶ Both lower and upper bounds : [155, 155]
- The number of intervals for m1 including 155 is equal to 8

Déclarative method:

• **Itemsets:** [Dzyuba et al., DAMI 2017]

Multisteps methods:

- **Itemsets:** [Boley et al., KDD 2011]
- Sequences: [Diop et al., ICDM 2018]
- Numerical data: [Giacometti et al., SDM 2018]



Theoretical proofs for the sampling distributions

5. Experimental Results



- cnrs
- Extending our approach to other interestingness measures like density
- Sampling patterns under constraints

- Incorporating our approach in interactive mining approaches
- Exploring the use of sampling in machine learning contexts