

The background features a complex, abstract design. It includes a network of thin, light-colored lines forming a web-like structure. Overlaid on this are various data visualizations: a grid of small, light-colored plus signs, a series of small, colorful dots (green, blue, yellow) connected by lines, and a large, semi-transparent white banner with a black border. The banner contains the main title in a bold, black, sans-serif font. The overall color palette is muted, with shades of brown, beige, and light gray, accented by the colors of the data points.

# **Session 6. Constraint-Based Sequential-Pattern Mining**

# Constraint-Based Sequential-Pattern Mining

---

- ❑ Share many similarities with constraint-based itemset mining
- ❑ **Anti-monotonic:** If  $S$  violates  $c$ , the super-sequences of  $S$  also violate  $c$ 
  - ❑  $\text{sum}(S.\text{price}) < 150; \min(S.\text{value}) > 10$
- ❑ **Monotonic:** If  $S$  satisfies  $c$ , the super-sequences of  $S$  also do so
  - ❑  $\text{element\_count}(S) > 5; S \supseteq \{\text{PC}, \text{digital\_camera}\}$
- ❑ **Data anti-monotonic:** If a sequence  $s_1$  with respect to  $S$  violates  $c_3$ ,  $s_1$  can be removed
  - ❑  $c_3: \text{sum}(S.\text{price}) \geq v$
- ❑ **Succinct:** Enforce constraint  $c$  by explicitly manipulating data
  - ❑  $S \supseteq \{\text{i-phone}, \text{MacAir}\}$
- ❑ **Convertible:** Projection based on the sorted value not sequence order
  - ❑  $\text{value\_avg}(S) < 25; \text{profit\_sum}(S) > 160$
  - ❑  $\text{max}(S)/\text{avg}(S) < 2; \text{median}(S) - \text{min}(S) > 5$


# Timing-Based Constraints in Seq.-Pattern Mining

---

- ❑ **Order constraint:** Some items must happen before the other
  - ❑  $\{\text{algebra, geometry}\} \rightarrow \{\text{calculus}\}$  (where “ $\rightarrow$ ” indicates ordering)
  - ❑ Anti-monotonic: Constraint-violating sub-patterns pruned
- ❑ **Min-gap/max-gap constraint:** Confines two elements in a pattern
  - ❑ E.g.,  $\text{mingap} = 1, \text{maxgap} = 4$
  - ❑ Succinct: Enforced directly during pattern growth
- ❑ **Max-span constraint:** Maximum allowed time difference between the 1<sup>st</sup> and the last elements in the pattern
  - ❑ E.g.,  $\text{maxspan}(S) = 60$  (days)
  - ❑ Succinct: Enforced directly when the 1<sup>st</sup> element is determined
- ❑ **Window size constraint:** Events in an element do not have to occur at the same time: Enforce max allowed time difference
  - ❑ E.g.,  $\text{window-size} = 2$ : Various ways to merge events into elements

# Episodes and Episode Pattern Mining

---

- Episodes and regular expressions: Alternative to seq. patterns
  - Serial episodes:  $A \rightarrow B$
  - Parallel episodes:  $A \mid B$   Indicating partial order relationships
  - Regular expressions:  $(A \mid B)C^*(D \rightarrow E)$
- Methods for episode pattern mining
  - Variations of Apriori/GSP-like algorithms
  - Projection-based pattern growth
    - $Q_1$ : Can you work out the details?
  - $Q_2$ : What are the differences between mining episodes and constraint-based pattern mining?

# Summary

---

- ❑ Concepts of Sequential Pattern Mining
- ❑ Sequential Pattern Mining Algorithms
  - ❑ **GSP** (Generalized Sequential Patterns)
  - ❑ Vertical Format-Based Mining: **SPADE**
  - ❑ Pattern-Growth Methods: **PrefixSpan**
- ❑ Mining Closed Sequential Patterns: **CloSpan**
- ❑ Constrain-Based Sequential Pattern Mining

# Recommended Readings

---

- ❑ M. N. Garofalakis, R. Rastogi, K. Shim: Mining Sequential Patterns with Regular Expression Constraints. IEEE Trans. Knowl. Data Eng. 14(3), 2002
- ❑ H. Mannila, H. Toivonen, and A. I. Verkamo, “Discovery of frequent episodes in event sequences”, Data Mining and Knowledge Discovery, 1997
- ❑ J. Pei, J. Han, B. Mortazavi-Asl, J. Wang, H. Pinto, Q. Chen, U. Dayal, and M.-C. Hsu, "Mining Sequential Patterns by Pattern-Growth: The PrefixSpan Approach", IEEE TKDE, 16(10), 2004
- ❑ J. Pei, J. Han, and W. Wang, "Constraint-based sequential pattern mining: the pattern-growth methods", J. Int. Inf. Sys., 28(2), 2007
- ❑ R. Srikant and R. Agrawal, “Mining sequential patterns: Generalizations and performance improvements”, EDBT’96
- ❑ X. Yan, J. Han, and R. Afshar, “CloSpan: Mining Closed Sequential Patterns in Large Datasets”, SDM'03
- ❑ M. Zaki, “SPADE: An Efficient Algorithm for Mining Frequent Sequences”, Machine Learning, 2001