

723A75 Advanced Data Mining TDDD41 Data Mining - Clustering and Association Analysis

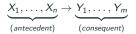
Lecture 9: Summary and Exercise

Johan Alenlöv IDA, Linköping University, Sweden

Outline

- Content
 - Summary
 - Exercise

- Goal: Given a transactional database find association rules on the form



with a user-specified minimum support and confidence.

- Support: The fraction of transactions that contains the full rule $X \cup Y$. $(p(X \cup Y))$
- Confidence: The fraction of transactions that contains X that also contains Y. (p(Y|X))
- Why? Help with decision making.
- Note that association is not causality.
- Two step solution:
 - 1. Generate all itemsets with a given minimum support.
 - 2. Generate all rules from these itemsets with minimum confidence.

- For generating frequent itemsets the following apriori property is important.
 - Every subset of a frequent itemset is frequent.
 - Alternatively every superset of an infrequent itemset is infrequent.
- Two algorithms for generating frequent itemsets
 - Apriori algorithm Using the apriori property to generate candidate sets that are tested.

Use the sets of length k to generate and test candidates of length k + 1.

FP Grow Construct an FP-tree and find the itemsets by looking at the conditional databases.

Constructs itemsets by building the chains with specific suffixes first.

- Given a frequent itemset *L* we wish to find a subset $X \subseteq L$ such that the rule $X \to L \setminus X$ has minimum confidence.
- Using the following property, if $X' \subseteq X$ then

$$\operatorname{Conf}(X \to L \setminus X) \ge \operatorname{Conf}(X' \to L \setminus X'),$$

we can reduce the number of sets to check.

• The algorithm goes over each subset (starting with maximal size) and then checking all subsets to find rules with minimum support.

- Other constraints can be added, such as minimum price, range of prices, sum of prices, etc.
- · Constraints can be,

Monotone If it is true for a set X then it is true for every superset X'. $(X \subseteq X')$ Antimonotone If it is true for a set X then it is true for every subset X'. $(X \supseteq X')$

Convertible Monotone If the items are sorted (in some way) then it is **monotone**. **Convertible Antimonotone** If the items are sorted (in some way) then it is

antimonotone.

Strongly convertible If it is both convertible monotone and convertible antimonotone.

Inconvertible Can't be converted.

 Run the Apriori algorithm on the following transactional database with minimum support equal to one transaction. Explain step by step the execution.

Tid	ltems
1	A, B, C
2	X, Y, Z
3	A, Y, C
4	X, B, Z

- Repeat with the constraint that the itemsets has to contain A. Make it clear when the constraint is used, don't just run the algorithm and consider the constraint at the end.
- Let the items A, B, C, X, Y, and Z, have the price of respectively
 -3, -2, -1, 1, 2, and 3 units. Repeat the exercise with the constraint: Find the frequent itemsets with range less than 3. Make it clear when the constraint is used, don't just run the algorithm and consider the constraint at the end.
- Repeat the exercises above with the FP grow algorithm
- Apply the rule generation algorithm to the frequent itemset XBZ on the database above in order to find association rules with confidence 0.5