# Mining Multiple-Level Associations

# Mining Multiple-Level Frequent Patterns

- ❑ Items often form hierarchies

    - ❑ Ex.: Dairyland 2% milk; Wonder wheat bread

- ❑ How to set min-support thresholds?

    - ❑ Uniform min-support across multiple levels (reasonable?)

    - ❑ Level-reduced min-support: Items at the lower level are expected to have lower support
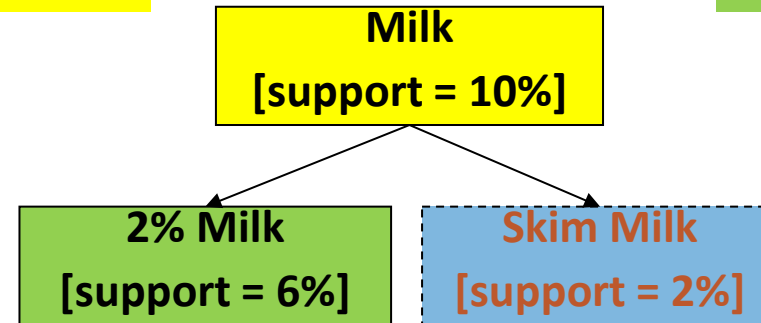
- ❑ Efficient mining: *Shared* multi-level mining

    - ❑ Use the lowest min-support to pass down the set of candidates

**Uniform support**

Level 1
min_sup = 5%

Level 2
min_sup = 5%

**Milk**
**[support = 10%]**

**2% Milk**
**[support = 6%]**

**Skim Milk**
**[support = 2%]**

**Reduced support**

Level 1
min_sup = 5%

Level 2
min_sup = 1%

# Redundancy Filtering at Mining Multi-Level Associations

❑ Multi-level association mining may generate many redundant rules

❑ Redundancy filtering:  Some rules may be redundant due to "ancestor" relationships between items

(Suppose the 2% milk sold is about ¼ of milk sold in gallons)

❑ milk $\Rightarrow$ wheat bread  [support = 8%, confidence = 70%]   (1)

❑ 2% milk $\Rightarrow$ wheat bread [support = 2%, confidence = 72%] (2)

❑ A rule is *redundant* if its support is close to the "expected" value, according to its "ancestor" rule, and it has a similar confidence as its "ancestor"

❑ Rule (1) is an ancestor of rule (2), which one to prune?

# Customized Min-Supports for Different Kinds of Items

❑ We have used the same min-support threshold for all the items or item sets to be mined in each association mining

❑ In reality, some items (e.g., diamond, watch, …) are valuable but less frequent

❑ It is necessary to have customized min-support settings for different kinds of items

❑ One Method: Use group-based "individualized" min-support

    ❑ E.g., {diamond, watch}: 0.05%; {bread, milk}: 5%; …

    ❑ How to mine such rules efficiently?

      ❑ Existing scalable mining algorithms can be easily extended to cover such cases

# Mining Multi-Dimensional Associations

- Single-dimensional rules (e.g., items are all in "product" dimension)

  - buys(X, "milk") $\Rightarrow$ buys(X, "bread")

- Multi-dimensional rules (i.e., items in $\geq$ 2 dimensions or predicates)

  - Inter-dimension association rules (*no repeated predicates*)

    - age(X, "18-25") $\wedge$ occupation(X, "student") $\Rightarrow$ buys(X, "coke")

  - Hybrid-dimension association rules (*repeated predicates*)

    - age(X, "18-25") $\wedge$ buys(X, "popcorn") $\Rightarrow$ buys(X, "coke")

- Attributes can be categorical or numerical

  - Categorical Attributes (e.g., *profession, product*: no ordering among values): Data cube for inter-dimension association

  - Quantitative Attributes: Numeric, implicit ordering among values—discretization, clustering, and gradient approaches

# Mining Quantitative Associations

- ❑ Mining associations with numerical attributes

  - ❑ Ex.: Numerical attributes: age and salary

- ❑ Methods

  - ❑ Static discretization based on predefined concept hierarchies

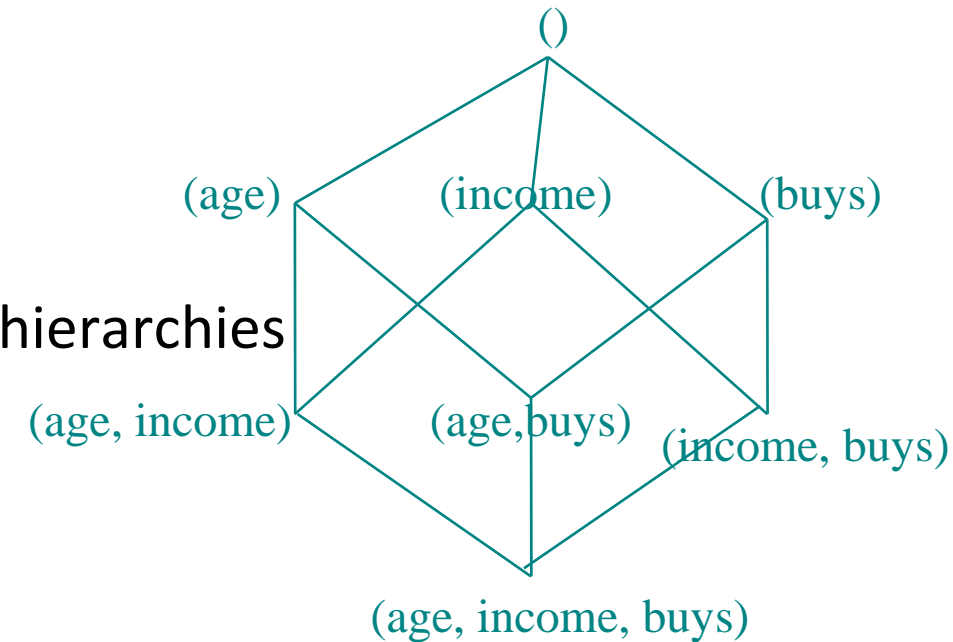    - ❑ Data cube-based aggregation

  - ❑ Dynamic discretization based on data distribution

  - ❑ Clustering: Distance-based association

    - ❑ First one-dimensional clustering, then association

  - ❑ Deviation analysis:

    - ❑ Gender = female $\Rightarrow$ Wage: mean=$7/hr (overall mean = $9)

()

(age)   (income)   (buys)

(age, income)   (age, buys)   (income, buys)

(age, income, buys)

# Mining Extraordinary Phenomena in Quantitative Association Mining

- ❑ Mining extraordinary (i.e., interesting) phenomena

  - ❑ Ex.: Gender = female ⇒ Wage: mean=$7/hr (overall mean = $9)

  - ❑ LHS: a subset of the population

  - ❑ RHS: an extraordinary behavior of this subset

- ❑ The rule is accepted only if a statistical test (e.g., Z-test) confirms the inference with high confidence

- ❑ Subrule: Highlights the extraordinary behavior of a subset of the population of the super rule

  - ❑ Ex.: (Gender = female) ^ (South = yes) ⇒ mean wage = $6.3/hr

- ❑ Rule condition can be categorical or numerical (quantitative rules)

  - ❑ Ex.: Education in [14-18] (yrs) ⇒ mean wage = $11.64/hr

- ❑ Efficient methods have been developed for mining such extraordinary rules (e.g., Aumann and Lindell@KDD'99)

Mining Negative Correlations

# Rare Patterns vs. Negative Patterns

❏ Rare patterns

  ❏ Very low support but interesting (e.g., buying Rolex watches)

  ❏ How to mine them? Setting individualized, group-based min-support thresholds for different groups of items

❏ Negative patterns

  ❏ Negatively correlated: Unlikely to happen together

  ❏ Ex.: Since it is unlikely that the same customer buys both a Ford Expedition (an SUV car) and a Ford Fusion (a hybrid car), buying a Ford Expedition and buying a Ford Fusion are likely negatively correlated patterns

  ❏ How to define negative patterns?

# Defining Negative Correlated Patterns

❑ A support-based definition

  ❑ If itemsets A and B are both frequent but rarely occur together, i.e.,
    sup(A U B) << sup (A) × sup(B)

  ❑ Then A and B are negatively correlated

  > Does this remind you the definition of *lift?*

❑ Is this a good definition for large transaction datasets?

❑ Ex.:  Suppose a store sold two needle packages A and B 100 times each, but only one transaction contained both A and B

  ❑ When there are in total 200 transactions, we have

    ❑ s(A U B) = 0.005, s(A) × s(B) = 0.25, s(A U B) << s(A) × s(B)

  ❑ But when there are $10^5$ transactions, we have

    ❑ s(A U B) = $1/10^5$, s(A) × s(B) = $1/10^3 × 1/10^3$, s(A U B) > s(A) × s(B)

  ❑ What is the problem?—Null transactions: The support-based definition is not null-invariant!

# Defining Negative Correlation: Need Null-Invariance in Definition

- ❑ A good definition on negative correlation should take care of the null-invariance problem

  - ❑ Whether two itemsets A and B are negatively correlated should not be influenced by the number of null-transactions

- ❑ A Kulczynski measure-based definition

  - ❑ If itemsets A and B are frequent but $(P(A|B) + P(B|A))/2 < \epsilon$, where $\epsilon$ is a negative pattern threshold, then A and B are negatively correlated

- ❑ For the same needle package problem:

  - ❑ No matter there are in total 200 or $10^5$ transactions

  - ❑ If $\epsilon = 0.01$, we have $(P(A|B) + P(B|A))/2 = (0.01 + 0.01)/2 < \epsilon$

# Mining Compressed Patterns

# Mining Compressed Patterns

| Pat-ID | Item-Sets | Support |
|--------|-----------|---------|
| P1 | {38,16,18,12} | 205227 |
| P2 | {38,16,18,12,17} | 205211 |
| P3 | {39,38,16,18,12,17} | 101758 |
| P4 | {39,16,18,12,17} | 161563 |
| P5 | {39,16,18,12} | 161576 |

❑ Closed patterns
  ❑ P1, P2, P3, P4, P5
  ❑ Emphasizes too much on support
  ❑ There is no compression
❑ Max-patterns
  ❑ P3: information loss
❑ Desired output (a good balance):
  ❑ P2, P3, P4

❑ Why mining compressed patterns?
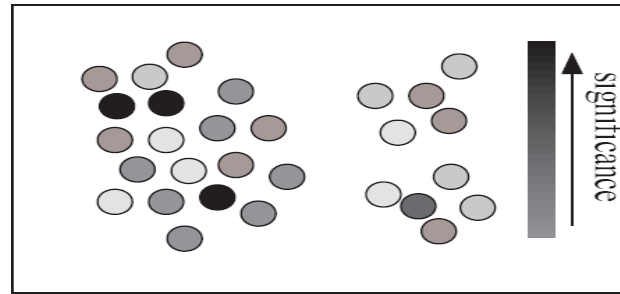  ❑ Too many scattered patterns but not so meaningful
❑ Pattern distance measure

$$Dist(P_1, P_2) = 1 - \frac{|T(P_1) \cap T(P_2)|}{|T(P_1) \cup T(P_2)|}$$
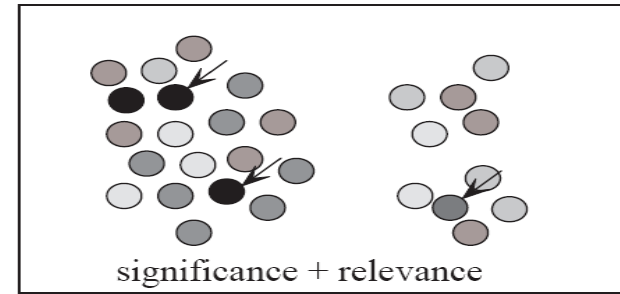
❑ δ-clustering: For each pattern P, find all patterns which can be expressed by P and whose distance to P is within δ (δ-cover)

❑ All patterns in the cluster can be represented by P

❑ Method for efficient, direct mining of compressed frequent patterns (e.g., D. Xin, J. Han, X. Yan, H. Cheng, "On Compressing Frequent Patterns", Knowledge and Data Engineering, 60:5-29, 2007)
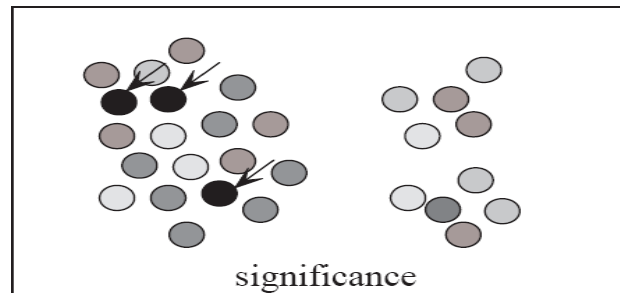
# Redundancy-Aware Top-k Patterns

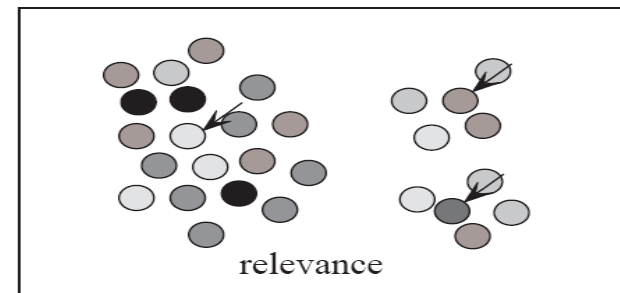❑ Desired patterns: high significance & low redundancy



(a) a set of patterns

(b) redundancy-aware top-$k$

(c) traditional top-$k$

(d) summarization

❑ Method: Use MMS (Maximal Marginal Significance) for measuring the combined significance of a pattern set

❑ Xin et al., Extracting Redundancy-Aware Top-K Patterns, KDD'06

# Summary

# Summary: Mining Diverse Patterns

❑   Efficient methods have been developed for mining various kinds of patterns

  ❑   Mining Multiple-Level Associations

  ❑   Mining Multi-Dimensional Associations

  ❑   Mining Quantitative Associations

  ❑   Mining Negative Correlations

  ❑   Mining Compressed and Redundancy-Aware Patterns

# Recommended Readings

- R. Srikant and R. Agrawal, "Mining generalized association rules", VLDB'95

- Y. Aumann and Y. Lindell, "A Statistical Theory for Quantitative Association Rules", KDD'99

- K. Wang, Y. He, J. Han, "Pushing Support Constraints Into Association Rules Mining", IEEE Trans. Knowledge and Data Eng. 15(3): 642-658, 2003

- D. Xin, J. Han, X. Yan and H. Cheng, "On Compressing Frequent Patterns", Knowledge and Data Engineering, 60(1): 5-29, 2007

- D. Xin, H. Cheng, X. Yan, and J. Han, "Extracting Redundancy-Aware Top-K Patterns", KDD'06

- J. Han, H. Cheng, D. Xin, and X. Yan, "Frequent Pattern Mining: Current Status and Future Directions", Data Mining and Knowledge Discovery, 15(1): 55-86, 2007