

What Is Pattern Discovery?

- What are patterns?
 - Patterns: A set of items, subsequences, or substructures that occur frequently together (or strongly correlated) in a data set
 - Patterns represent intrinsic and important properties of datasets
- Pattern discovery: Uncovering patterns from massive data sets
- Motivation examples:
 - What products were often purchased together?
 - What are the subsequent purchases after buying an iPad?
 - What code segments likely contain copy-and-paste bugs?
 - What word sequences likely form phrases in this corpus?

4#

Pattern Discovery: Why Is It Important?

- Finding inherent regularities in a data set
- Foundation for many essential data mining tasks
 - Association, correlation, and causality analysis
 - Mining sequential, structural (e.g., sub-graph) patterns
 - Pattern analysis in spatiotemporal, multimedia, time-series, and stream data
 - Classification: Discriminative pattern-based analysis
 - Cluster analysis: Pattern-based subspace clustering
- Broad applications
 - Market basket analysis, cross-marketing, catalog design, sale campaign analysis, Web log analysis, biological sequence analysis

4#

Basic Concepts: k-Itemsets and Their Supports

- Itemset: A set of one or more items
- k-itemset: $X = \{x_1, \dots, x_k\}$
- Ex. {Beer, Nuts, Diaper} is a 3-itemset
- (absolute) support (count) of X, $\text{sup}\{X\}$: Frequency or the number of occurrences of an itemset X
 - Ex. $\text{sup}\{\text{Beer}\} = 3$
 - Ex. $\text{sup}\{\text{Diaper}\} = 4$
 - Ex. $\text{sup}\{\text{Beer, Diaper}\} = 3$
 - Ex. $\text{sup}\{\text{Beer, Eggs}\} = 1$
- (relative) support, $s\{X\}$: The fraction of transactions that contains X (i.e., the probability that a transaction contains X)
 - Ex. $s\{\text{Beer}\} = 3/5 = 60\%$
 - Ex. $s\{\text{Diaper}\} = 4/5 = 80\%$
 - Ex. $s\{\text{Beer, Eggs}\} = 1/5 = 20\%$

7 of 55

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

- (relative) support, $s\{X\}$: The fraction of transactions that contains X (i.e., the probability that a transaction contains X)
 - Ex. $s\{\text{Beer}\} = 3/5 = 60\%$
 - Ex. $s\{\text{Diaper}\} = 4/5 = 80\%$
 - Ex. $s\{\text{Beer, Eggs}\} = 1/5 = 20\%$

(threshold) 70% / 80%, 100%
↳

Basic Concepts: Frequent Itemsets (Patterns)

- An itemset (or a pattern) X is frequent if the support of X is no less than a minsup threshold σ ↳ threshold σ
- Let $\sigma = 50\%$ (σ : minsup threshold)
For the given 5-transaction dataset
 - All the frequent 1-itemsets:
 - Beer: $3/5 (60\%)$; Nuts: $3/5 (60\%)$
 - Diaper: $4/5 (80\%)$; Eggs: $3/5 (60\%)$
 - All the frequent 2-itemsets:
 - {Beer, Diaper}: $3/5 (60\%)$
 - All the frequent 3-itemsets?
 - None
- Why do these itemsets (shown on the left) form the complete set of frequent k-itemsets (patterns) for any k ?
- Observation: We may need an efficient method to mine a complete set of frequent patterns
↳ coffee : $2/5 (40\%)$ < minsup threshold ↳

From Frequent Itemsets to Association Rules

- Comparing with itemsets, rules can be more telling
 - Ex. $\text{Diaper} \rightarrow \text{Beer}$ \rightarrow Diaper \rightarrow Beer \rightarrow Beer \rightarrow Beer
 - Buying diapers may likely lead to buying beers*
- How strong is this rule? (support, confidence)
 - Measuring association rules: $X \rightarrow Y$ (s, c)
 - Both X and Y are itemsets
 - Support**, s : The probability that a transaction contains $X \cup Y$ \rightarrow $\frac{\text{Number of transactions containing } X \cup Y}{\text{Total number of transactions}}$
 - Ex. $s\{\text{Diaper}, \text{Beer}\} = 3/5 = 0.6$ (i.e., 60%)
 - Confidence**, c : The conditional probability that a transaction containing X also contains Y
 - Calculation: $c = \text{sup}(X \cup Y) / \text{sup}(X)$
 - Note: $X \cup Y$: the union of two itemsets
 - Ex. $c = \text{sup}\{\text{Diaper}, \text{Beer}\} / \text{sup}\{\text{Diaper}\} = \frac{3}{4} = 75\%$

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

Containing both beer {Beer} \cup {Diaper} = {Beer, Diaper}

Note: $X \cup Y$: the union of two itemsets
The set contains both X and Y

9 of 55

Mining Frequent Itemsets and Association Rules

- Association rule mining** \rightarrow Mining itemsets
 - Given two thresholds: minsup , minconf
 - Find all of the rules, $X \rightarrow Y$ (s, c)
 - such that, $s \geq \text{minsup}$ and $c \geq \text{minconf}$
 - Let $\text{minsup} = 50\%$ \rightarrow itemset $\{Beer, Diaper\}$ has support 3
 - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
 - Freq. 2-itemsets: $\{Beer, Diaper\}: 3$
 - Let $\text{minconf} = 50\%$ \rightarrow itemset $\{Beer, Diaper\}$ has confidence 100%
 - $Beer \rightarrow Diaper$ (60%, 100%) $\rightarrow \frac{\text{sup}(Beer \cup Diaper)}{\text{sup}(Beer)}$
 - $Diaper \rightarrow Beer$ (10%, 75%) $\rightarrow \frac{\text{sup}(Beer \cup Diaper)}{\text{sup}(Diaper)}$

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

Observations:

- Mining association rules and mining frequent patterns are very close problems
- Scalable methods are needed for mining large datasets

#

Chapter 6: Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

- Basic Concepts
- Efficient Pattern Mining Methods \rightarrow
- Pattern Evaluation
- Summary

#

Efficient Pattern Mining Methods

- The Downward Closure Property of Frequent Patterns
- The **Apriori Algorithm**
- Extensions or Improvements of Apriori
- Mining Frequent Patterns by Exploring Vertical Data Format
- FPGrowth: A Frequent Pattern-Growth Approach
- Mining Closed Patterns

#

Apriori Pruning and Scalable Mining Methods

- Apriori pruning principle: If there is any itemset which is infrequent, its superset should not even be generated! (Agrawal & Srikant @VLDB'94, Mannila, et al. @ KDD' 94)
- Scalable mining Methods: Three major approaches
 - Level-wise, join-based approach: Apriori (Agrawal & Srikant@VLDB'94)
 - Vertical data format approach: Eclat (Zaki, Parthasarathy, Ogihara, Li @KDD'97)
 - Frequent pattern projection and growth: FPgrowth (Han, Pei, Yin @SIGMOD'00)

ମୁଣ୍ଡିଲେ
ପରିବାହନ ପାଇଁ
କ୍ଷେତ୍ରକୁ ଦେଖିବାରେ

Apriori: A Candidate Generation & Test Approach

- Outline of Apriori (level-wise, candidate generation and test)
 - Initially, scan DB once to get frequent 1-itemset
 - **Repeat**
 - Generate length-(k+1) candidate itemsets from length-k frequent itemsets
 - Test the candidates against DB to find frequent (k+1)-itemsets
 - Set k := k +1
 - Until no frequent or candidate set can be generated
 - Return all the frequent itemsets derived

The Apriori Algorithm—An Example

