CSE 5243 INTRO. TO DATA MINING

Mining Frequent Patterns and Associations: Basic Concepts Yu Su, CSE@The Ohio State University

Slides adapted from UIUC CS412 by Prof. Jiawei Han and OSU CSE5243 by Prof. Huan Sun

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





- Efficient Pattern Mining Methods
- Pattern Evaluation



Pattern Discovery: Basic Concepts

What Is Pattern Discovery? Why Is It Important?

Basic Concepts: Frequent Patterns and Association Rules

Compressed Representation: Closed Patterns and Max-Patterns

What Is Pattern Discovery?

- Motivating 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?

What Is Pattern Discovery?

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?

□ 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

What Is Pattern Discovery?

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?

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

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

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

□ **Itemset**: A set of one or more items

- Itemset: A set of one or more items
- □ **k-itemset**: $X = \{x_1, ..., x_k\}$
 - **Ex.** {Beer, Nuts, Diaper} is a 3-itemset

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

- Itemset: A set of one or more items
- □ **k-itemset**: $X = \{x_1, ..., x_k\}$
 - Ex. {Beer, Nuts, Diaper} is a 3-itemset
- (absolute) support (count) of X, sup{X}: Frequency or the number of occurrences of an itemset X
 - **Ex.** sup{Beer} = 3
 - Ex. sup{Diaper} = 4
 - Ex. sup{Beer, Diaper} = 3
 - Ex. sup{Beer, Eggs} = 1

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

- Itemset: A set of one or more items
- □ **k-itemset**: $X = \{x_1, ..., x_k\}$
 - Ex. {Beer, Nuts, Diaper} is a 3-itemset
- (absolute) support (count) of X, sup{X}: Frequency or the number of occurrences of an itemset X
 - Ex. sup{Beer} = 3
 - Ex. sup{Diaper} = 4
 - Ex. sup{Beer, Diaper} = 3
 - Ex. sup{Beer, Eggs} = 1

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{Beer} = 3/5 = 60%
 - Ex. $s{Diaper} = 4/5 = 80\%$
 - **Ex.** s{Beer, Eggs} = 1/5 = 20%

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

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
- Let σ = 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

Ti	id	Items bought
1	.0	Beer, Nuts, Diaper
2	20	Beer, Coffee, Diaper
3	0	Beer, Diaper, Eggs
4	0	Nuts, Eggs, Milk
5	50	Nuts, Coffee, Diaper, Eggs, Milk

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
- Let σ = 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

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

- Do these itemsets (shown on the left) form the complete set of frequent kitemsets (patterns) for any k?
- Observation: We may need an efficient method to mine a complete set of frequent patterns

- Comparing with itemsets, rules can be more telling
 - \square Ex. Diaper \rightarrow Beer
 - Buying diapers may likely lead to buying beers

- Ex. Diaper → 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

- Ex. Diaper → 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
 - Ex. s{Diaper, Beer} = 3/5 = 0.6 (i.e., 60%)

	Tid	Items bought
	10	Beer, Nuts, Diaper
	20	Beer, Coffee, Diaper
	30	Beer, Diaper, Eggs
	40	Nuts, Eggs, Milk
1	50	Nuts, Coffee, Diaper, Eggs, Milk

- Ex. Diaper → 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
 - Ex. s{Diaper, Beer} = 3/5 = 0.6 (i.e., 60%)
 - Confidence, c: The conditional probability that a transaction containing X also contains Y
 - Calculation: $c = sup(X \cup Y) / sup(X)$
 - Ex. c = sup{Diaper, Beer}/sup{Diaper} = ³/₄ = 0.75

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



Association rule mining

- Given two thresholds: *minsup, minconf*
- **The Find all of the rules**, $X \rightarrow Y$ (s, c)
 - such that, $s \ge minsup$ and $c \ge minconf$

Association rule mining

- Given two thresholds: *minsup, minconf*
- **The Find all of the rules**, $X \rightarrow Y$ (s, c)
 - such that, $s \ge minsup$ and $c \ge minconf$
- $\Box \quad \text{Let } minsup = 50\%$
 - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
 - □ Freq. 2-itemsets: {Beer, Diaper}: 3

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

Association rule mining

- Given two thresholds: *minsup, minconf*
- **The Find all of the rules**, $X \rightarrow Y$ (s, c)
 - such that, $s \ge minsup$ and $c \ge minconf$
- $\Box \quad \text{Let } minsup = 50\%$
 - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
 - □ Freq. 2-itemsets: {Beer, Diaper}: 3
- $\Box \quad \text{Let minconf} = 50\%$
 - $\square \quad Beer \rightarrow Diaper \quad (60\%, 100\%)$
 - $\Box \quad \text{Diaper} \rightarrow \text{Beer} \quad (60\%, 75\%)$

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

Association rule mining

- Given two thresholds: *minsup, minconf*
- **The Find all of the rules**, $X \rightarrow Y$ (s, c)
 - such that, $s \ge minsup$ and $c \ge minconf$
- $\Box \quad \text{Let } minsup = 50\%$
 - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
 - □ Freq. 2-itemsets: {Beer, Diaper}: 3
- $\Box \quad \text{Let minconf} = 50\%$
 - $\square \quad Beer \rightarrow Diaper (60\%, 100\%)$
 - $\Box \quad \text{Diaper} \rightarrow \text{Beer} \quad (60\%, 75\%)$

(Q: Are these all rules?)

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

Association rule mining

- Given two thresholds: *minsup, minconf*
- **•** Find all of the rules, $X \rightarrow Y$ (s, c)
 - such that, $s \ge minsup$ and $c \ge minconf$
- \Box Let minsup = 50%
 - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
 - □ Freq. 2-itemsets: {Beer, Diaper}: 3
- $\Box \quad \text{Let minconf} = 50\%$
 - $\square \quad Beer \rightarrow Diaper \quad (60\%, 100\%)$
 - □ Diaper → Beer (60%, 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

Observations:

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

Association Rule Mining: two-step process

In general, association rule mining can be viewed as a two-step process:

- 1. Find all frequent itemsets: By definition, each of these itemsets will occur at least as frequently as a predetermined minimum support count, *min_sup*.
- 2. Generate strong association rules from the frequent itemsets: By definition, these rules must satisfy minimum support and minimum confidence.

Because the second step is much less costly than the first, the overall performance of mining association rules is determined by the first step.

Generating Association Rules from Frequent Patterns

□ Recall that: $confidence(A \Rightarrow B) = P(B|A) = \frac{support_count(A \cup B)}{support_count(A)}$

Once we mined frequent patterns, association rules can be generated as follows:

- For each frequent itemset l, generate all nonempty subsets of l.
- For every nonempty subset s of l, output the rule " $s \Rightarrow (l-s)$ " if $\frac{support_count(l)}{support_count(s)} \ge min_conf$, where min_conf is the minimum confidence threshold.

Generating Association Rules from Frequent Patterns

□ Recall that:

 $confidence(A \Rightarrow B) = P(B|A) = \frac{support_count(A \cup B)}{support_count(A)}$

Once we mined frequent patterns, association rules can be generated as follows:

- For each frequent itemset l, generate all nonempty subsets of l.
- For every nonempty subset s of l, output the rule " $s \Rightarrow (l s)$ " if $\frac{support_count(l)}{support_count(s)} \ge min_conf$, where min_conf is the minimum confidence threshold.

Because l is a frequent itemset, each rule automatically satisfies the minimum support requirement.

Example: Generating Association Rules

Generating association rules. Let's try an example based on the transactional data for *AllElectronics* shown in Table 6.1. The data contain frequent itemset $X = \{I1, I2, I5\}$. What are the association rules that can be generated from X? The nonempty subsets of X are $\{I1, I2\}, \{I1, I5\}, \{I2, I5\}, \{I1\}, \{I2\},$ and $\{I5\}$. The resulting association rules are as shown below, each listed with its confidence:

 $\{I1, I2\} \Rightarrow I5, \\ \{I1, I5\} \Rightarrow I2, \\ \{I2, I5\} \Rightarrow I1, \\ I1 \Rightarrow \{I2, I5\}, \\ I2 \Rightarrow \{I1, I5\}, \\ I5 \Rightarrow \{I1, I2\},$

confidence = 2/4 = 50% confidence = 2/2 = 100% confidence = 2/2 = 100% confidence = 2/6 = 33% confidence = 2/7 = 29%confidence = 2/2 = 100%

If minimum confidence threshold: 70%, what will be output?

Example from Chapter 6

Challenge: There Are Too Many Frequent Patterns!

□ A long pattern contains a combinatorial number of sub-patterns

- \Box How many frequent itemsets does the following TDB₁ contain?
 - **TDB**_{1:} $T_1: \{a_1, ..., a_{50}\}; T_2: \{a_1, ..., a_{100}\}$
 - Assuming (absolute) minsup = 1
 - Let's give it a try...
 - 1-itemsets: $\{a_1\}$: 2, $\{a_2\}$: 2, ..., $\{a_{50}\}$: 2, $\{a_{51}\}$: 1, ..., $\{a_{100}\}$: 1,

2-itemsets: $\{a_1, a_2\}$: 2, ..., $\{a_1, a_{50}\}$: 2, $\{a_1, a_{51}\}$: 1 ..., ..., $\{a_{99}, a_{100}\}$: 1,

99-itemsets: {a₁, a₂, ..., a₉₉}: 1, ..., {a₂, a₃, ..., a₁₀₀}: 1
100-itemset: {a₁, a₂, ..., a₁₀₀}: 1

Challenge: There Are Too Many Frequent Patterns!

- A long pattern contains a combinatorial number of sub-patterns
- □ How many frequent itemsets does the following TDB₁ contain?
 - **TDB**_{1:} $T_1: \{a_1, ..., a_{50}\}; T_2: \{a_1, ..., a_{100}\}$
 - Assuming (absolute) minsup = 1
 - Let's give it a try...
 - 1-itemsets: $\{a_1\}$: 2, $\{a_2\}$: 2, ..., $\{a_{50}\}$: 2, $\{a_{51}\}$: 1, ..., $\{a_{100}\}$: 1,
 - 2-itemsets: $\{a_1, a_2\}$: 2, ..., $\{a_1, a_{50}\}$: 2, $\{a_1, a_{51}\}$: 1 ..., ..., $\{a_{99}, a_{100}\}$: 1,

••••, ••••, ••••, •••

99-itemsets: {a₁, a₂, ..., a₉₉}: 1, ..., {a₂, a₃, ..., a₁₀₀}: 1 100-itemset: {a₁, a₂, ..., a₁₀₀}: 1

□ The total number of frequent itemsets:

$$\binom{100}{1} + \binom{100}{2} + \binom{100}{3} + \dots + \binom{100}{100} = 2^{100} - 1$$

Too huge a set for any one to compute or store!

Expressing Patterns in Compressed Form: Closed Patterns

- □ How to handle such a challenge?
- Solution 1: Closed patterns: A pattern (itemset) X is closed if X is frequent, and there exists no super-pattern Y Color X, with the same support as X

Expressing Patterns in Compressed Form: Closed Patterns

- □ How to handle such a challenge?
- Solution 1: Closed patterns: A pattern (itemset) X is closed if X is frequent, and there exists no super-pattern Y Construct X, with the same support as X
 - **Let Transaction DB TDB**₁: $T_1: \{a_1, ..., a_{50}\}; T_2: \{a_1, ..., a_{100}\}$
 - Suppose minsup = 1. How many closed patterns does TDB₁ contain?
 - Two: P_1 : "{ a_1 , ..., a_{50} }: 2"; P_2 : "{ a_1 , ..., a_{100} }: 1"

Why?

Expressing Patterns in Compressed Form: Closed Patterns

- □ How to handle such a challenge?
- Solution 1: Closed patterns: A pattern (itemset) X is closed if X is frequent, and there exists no super-pattern Y I X, with the same support as X
 - **Let Transaction DB TDB**₁: $T_1: \{a_1, ..., a_{50}\}; T_2: \{a_1, ..., a_{100}\}$
 - Suppose minsup = 1. How many closed patterns does TDB₁ contain?

• Two:
$$P_1$$
: "{ a_1 , ..., a_{50} }: 2"; P_2 : "{ a_1 , ..., a_{100} }: 1"

- Closed pattern is a lossless compression of frequent patterns
 - Reduces the # of patterns but does not lose the support information!
 - You will still be able to say: " $\{a_2, ..., a_{40}\}$: 2", " $\{a_5, a_{51}\}$: 1"

Expressing Patterns in Compressed Form: Max-Patterns

Solution 2: Max-patterns: A pattern X is a max-pattern if X is frequent and there exists no frequent super-pattern Y > X

Expressing Patterns in Compressed Form: Max-Patterns

- Solution 2: Max-patterns: A pattern X is a max-pattern if X is frequent and there exists no frequent super-pattern Y > X
- Difference with closed-patterns?
 - Do not care about the real support of the sub-patterns of a max-pattern
 - **Let Transaction DB TDB**₁: $T_1: \{a_1, ..., a_{50}\}; T_2: \{a_1, ..., a_{100}\}$
 - Suppose minsup = 1. How many max-patterns does TDB₁ contain?
 - One: P: "{a₁, ..., a₁₀₀}: 1"

Why?

Expressing Patterns in Compressed Form: Max-Patterns

- Solution 2: Max-patterns: A pattern X is a max-pattern if X is frequent and there exists no frequent super-pattern Y > X
- Difference with closed-patterns?
 - Do not care about the real support of the sub-patterns of a max-pattern
 - **Let Transaction DB TDB**₁: $T_1: \{a_1, ..., a_{50}\}; T_2: \{a_1, ..., a_{100}\}$
 - Suppose minsup = 1. How many max-patterns does TDB₁ contain?
 - One: P: "{a₁, ..., a₁₀₀}: 1"
- □ Max-pattern is a lossy compression!
 - We only know $\{a_1, ..., a_{40}\}$ is frequent
 - **D** But we do not know the real support of $\{a_1, \ldots, a_{40}\}, \ldots, any more \}$
 - Thus in many applications, closed-patterns are more desirable than max-patterns


Closed and maximal frequent itemsets. Suppose that a transaction database has only two transactions: $\{\langle a_1, a_2, \ldots, a_{100} \rangle; \langle a_1, a_2, \ldots, a_{50} \rangle\}$. Let the minimum support count threshold be $min_sup = 1$. We find two closed frequent itemsets and their support counts, that is, $C = \{\{a_1, a_2, \ldots, a_{100}\} : 1; \{a_1, a_2, \ldots, a_{50}\} : 2\}$. There is only one maximal frequent itemset: $\mathcal{M} = \{\{a_1, a_2, \ldots, a_{100}\} : 1\}$. Notice that we cannot include $\{a_1, a_2, \ldots, a_{50}\}$ as a maximal frequent itemset because it has a frequent super-set, $\{a_1, a_2, \ldots, a_{100}\}$. Compare this to the above, where we determined that there are $2^{100} - 1$ frequent itemsets, which is too huge a set to be enumerated!

{all frequent patterns} \supseteq {closed frequent patterns} \supseteq {max frequent patterns}



Closed and maximal frequent itemsets. Suppose that a transaction database has only two transactions: $\{\langle a_1, a_2, \ldots, a_{100} \rangle; \langle a_1, a_2, \ldots, a_{50} \rangle\}$. Let the minimum support count threshold be $min_sup = 1$. We find two closed frequent itemsets and their support counts, that is, $C = \{\{a_1, a_2, \ldots, a_{100}\} : 1; \{a_1, a_2, \ldots, a_{50}\} : 2\}$. There is only one maximal frequent itemset: $\mathcal{M} = \{\{a_1, a_2, \ldots, a_{100}\} : 1\}$. Notice that we cannot include $\{a_1, a_2, \ldots, a_{50}\}$ as a maximal frequent itemset because it has a frequent super-set, $\{a_1, a_2, \ldots, a_{100}\}$. Compare this to the above, where we determined that there are $2^{100}-1$ frequent itemsets, which is too huge a set to be enumerated!

The set of closed-patterns contains complete information regarding the frequent itemsets.

Quiz

□ Given closed frequent itemsets:

$$C = \{ \{a1, a2, ..., a100\}: 1; \{a1, a2, ..., a50\}: 2 \}$$

Is {a2, a45} frequent? Can we know its support?

Quiz (Cont'd)

□ Given maximal frequent itemset:

What is the support of {a8, a55}?

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

Basic Concepts

🗌 Efficient Pattern Mining Methods 🧏

The Apriori Algorithm

Application in Classification

Pattern Evaluation



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

The Downward Closure Property of Frequent Patterns

- Observation: From $TDB_{1:}T_1: \{a_1, ..., a_{50}\}; T_2: \{a_1, ..., a_{100}\}$
 - We get a frequent itemset: {a₁, ..., a₅₀}
 - Also, its subsets are all frequent: {a₁}, {a₂}, ..., {a₅₀}, {a₁, a₂}, ..., {a₁, ..., a₄₉}, ...
 - There must be some hidden relationships among frequent patterns!

The Downward Closure Property of Frequent Patterns

- Observation: From $TDB_{1:}T_1$: { $a_1, ..., a_{50}$ }; T_2 : { $a_1, ..., a_{100}$ }
 - We get a frequent itemset: {a₁, ..., a₅₀}
 - Also, its subsets are all frequent: {a₁}, {a₂}, ..., {a₅₀}, {a₁, a₂}, ..., {a₁, ..., a₄₉}, ...
 - There must be some hidden relationships among frequent patterns!
- □ The downward closure (also called "Apriori") property of frequent patterns
 - If {beer, diaper, nuts} is frequent, so is {beer, diaper}
 - Every transaction containing {beer, diaper, nuts} also contains {beer, diaper}
 - Apriori: Any subset of a frequent itemset must be frequent

A sharp knife for pruning!

The Downward Closure Property of Frequent Patterns

- Observation: From $TDB_{1:}T_1$: { $a_1, ..., a_{50}$ }; T_2 : { $a_1, ..., a_{100}$ }
 - We get a frequent itemset: {a₁, ..., a₅₀}
 - Also, its subsets are all frequent: {a₁}, {a₂}, ..., {a₅₀}, {a₁, a₂}, ..., {a₁, ..., a₄₉}, ...
 - There must be some hidden relationships among frequent patterns!
- □ The downward closure (also called "Apriori") property of frequent patterns
 - If {beer, diaper, nuts} is frequent, so is {beer, diaper}
 - Every transaction containing {beer, diaper, nuts} also contains {beer, diaper}
 - Apriori: Any subset of a frequent itemset must be frequent
- Efficient mining methodology

- A sharp knife for pruning!
- If any subset of an itemset S is infrequent, then there is no chance for S to be frequent—why do we even have to consider S ?!

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)

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 (Pseudo-Code)

```
C<sub>k</sub>: Candidate itemset of size k
F_k: Frequent itemset of size k
k := 1;
F_k := \{ \text{frequent items} \}; // \text{frequent 1-itemset} 
While (F_k != \emptyset) do \{ // \text{when } F_k \text{ is non-empty} \}
    C_{k+1} := candidates generated from F_k; // candidate generation
    Derive F_{k+1} by counting candidates in C_{k+1} with respect to TDB at minsup;
    k := k + 1
return \cup_{\nu} F_{\nu}
                           // return F_{k} generated at each level
```

The Apriori Algorithm—An Example



The Apriori Algorithm—An Example



Apriori: Implementation Tricks

- □ How to generate candidates?
 - **Step 1:** self-joining F_k
 - □ Step 2: pruning

Apriori: Implementation Tricks

- How to generate candidates?
 - **D** Step 1: self-joining F_k
 - Step 2: pruning
- Example of candidate-generation
 - $\blacksquare F_3 = \{abc, abd, acd, ace, bcd\}$
 - **Self-joining:** $F_3 * F_3$
 - abcd from abc and abd
 - acde from acd and ace



Apriori: Implementation Tricks

- How to generate candidates?
 - **D** Step 1: self-joining F_k
 - Step 2: pruning
- Example of candidate-generation
 - $\blacksquare F_3 = \{abc, abd, acd, ace, bcd\}$
 - **Self-joining:** $F_3 * F_3$
 - abcd from abc and abd
 - acde from acd and ace
 - Pruning:
 - acde is removed because ade is not in F_3
 - $\Box C_4 = \{ abcd \}$



Candidate Generation: An SQL Implementation



Apriori Adv/Disadv

Advantages:

- Uses large itemset property
- Easily parallelized
- Easy to implement

Disadvantages:

Assumes transaction database is memory resident

Requires up to m database scans

Classification based on Association Rules (CBA)

□ Why?

- Can effectively uncover the correlation structure in data
- AR are typically quite scalable in practice
- Rules are often very intuitive
 - Hence classifier built on intuitive rules is easier to interpret
- □ When to use?
 - On large dynamic datasets where class labels are available and the correlation structure is unknown.
 - Multi-class categorization problems
 - E.g. Web/Text Categorization, Network Intrusion Detection

Classification based on Association Rules (CBA)

Input

- <feature vector> <class label(s)>
- \Box <feature vector> = w1,...,wN
- class label(s)> = c1,...,cM

□ Run AR with minsup and minconf

- Prune rules of form
 - w1 → w2, [w1,c2] → c3 etc.

Keep only rules satisfying the constraints:

• $W \rightarrow C$ (Left: only composed of w1,...wN and Right: only composed of c1,...cM)

e.g., text categorization

CBA: Text Categorization (cont.)

Order remaining rules

- By confidence
 - **100%**
 - R1: W1 \rightarrow C1 (support 40%)
 - **R2:** W4 \rightarrow C2 (support 60%)

95%

- **R3:** W3 \rightarrow C2 (support 30%)
- **R4:** W5 \rightarrow C4 (support 70%)

And within each confidence level by support

Ordering R2, R1, R4, R3

Classification based Association

CBA: Text Categorization (cont.)

- Take training data and evaluate the predictive ability of each rule, prune rules that are subsumed by superior rules
 - **T1: W1 W5 C1,C4**
 - **T2:** W2 W4 C2
 - T3: W3 W4 C2
 - **T4: W5 W8 C4**

- Note: only a subset of transactions in training data
- Rule R3 would be pruned in this example if it is always subsumed by Rule R2 R3: $W3 \rightarrow C2$ R2: $W4 \rightarrow C2$ Why?

CBA: Text Categorization (cont.)

- Take training data and evaluate the predictive ability of each rule, prune rules that are subsumed by superior rules
 - □ T1: W1 W5 C1,C4
 - T2: W2 W4 C2
 - T3: W3 W4 C2
 - **T4: W5 W8 C4**

- Note: only a subset of transactions
- in training data

Rule R3 would be pruned in this example if it is always subsumed by Rule R2 {T3} is predictable by R3: W3 \rightarrow C2 {T2, T3} is predictable by R2: W4 \rightarrow C2 R3 is subsumed by R2, and will therefore be pruned.

Formal Concepts of Model

□ Given two rules r_i and r_i , define: $r_i \succ r_i$ if

The confidence of r_i is greater than that of r_i , or

Their confidences are the same, but the support of r_i is greater than that of r_j , or Both the confidences and supports are the same, but r_i is generated earlier than r_j .

Our classifier model is of the following format:

<r_1, r_2, ..., r_n, default_class>, where $r_i \in R$, $r_a \succ r_b$ if b > a

□ Other models possible

Sort by length of antecedent

Using the CBA model to classify

- For a new transaction
 - □ W1, W3, W5
 - Pick the k-most confident rules that apply (using the precedence ordering established in the baseline model)
 - The resulting classes are the predictions for this transaction
 - If k = 1 you would pick ?
 - If k = 2 you would pick ?

Conf: 100%
R1: W1 → C1 (support 40%)
R2: W4 → C2 (support 60%)
Conf: 95%
R3: W3 → C2 (support 30%)
R4: W5 → C4 (support 70%)

Using the CBA model to classify

- For a new transaction
 W1, W3, W5
 - Pick the k-most confident rules that apply (using the precedence ordering established in the baseline model)
 - The resulting classes are the predictions for this transaction
 - If k = 1 you would pick C1
 - If k = 2 you would pick C1, C4 (multi-class)
 - If W9, W10 (not covered by any rule), you would pick C2 (assuming it's the default, most dominant class)
 - Accuracy measurements as before (Classification Error)

CBA: Procedural Steps

Preprocessing, Training and Testing data split

- □ Compute AR on Training data
 - Keep only rules of form $X \rightarrow C$
 - C is class label itemset and X is feature itemset
- Order AR
 - According to confidence
 - According to support (at each confidence level)
- Prune away rules that lack sufficient predictive ability on training data (starting top-down)
 Rule subsumption
- □ For data that is not predictable, pick most dominant class as default class
- Test on testing data and report accuracy

Apriori: Improvements and Alternatives

Reduce passes of transaction database scans

Partitioning (e.g., Savasere, et al., 1995)

Dynamic itemset counting (Brin, et al., 1997)

- Shrink the number of candidates
 - Hashing (e.g., DHP: Park, et al., 1995)
 - Pruning by support lower bounding (e.g., Bayardo 1998)
 - Sampling (e.g., Toivonen, 1996)
- Exploring special data structures
 - Tree projection (Agarwal, et al., 2001)
 - H-miner (Pei, et al., 2001)
 - Hypecube decomposition (e.g., LCM: Uno, et al., 2004)





<1> Partitioning: Scan Database Only Twice

Theorem: Any itemset that is potentially frequent in TDB must be frequent in at least one of the partitions of TDB

Why?

<1> Partitioning: Scan Database Only Twice

Theorem: Any itemset that is potentially frequent in TDB must be frequent in at least one of the partitions of TDB



Proof by contradiction

<1> Partitioning: Scan Database Only Twice

Theorem: Any itemset that is potentially frequent in TDB must be frequent in at least one of the partitions of TDB



- Method: Scan DB twice (A. Savasere, E. Omiecinski and S. Navathe, VLDB'95)
 - Scan 1: Partition database so that each partition can fit in main memory
 - Mine local frequent patterns in this partition
 - Scan 2: Consolidate global frequent patterns
 - Find global frequent itemset candidates (those frequent in at least one partition)
 - Find the true frequency of those candidates, by scanning TDB_i one more time

<2> Direct Hashing and Pruning (DHP):

- Reduce candidate number: (J. Park, M. Chen, and P. Yu, SIGMOD'95)
- \Box Hashing: Different itemsets may have the same hash value: v = hash(itemset)
- □ 1st scan: When counting the 1-itemset, hash 2-itemset to calculate the bucket count
- Observation: A k-itemset cannot be frequent if its corresponding hashing bucket count is below the minsup threshold
- Example: At the 1st scan of TDB, count 1-itemset, and
 - Hash 2-itemsets in each transaction to its bucket
 - {ab, ad, ce}
 - {bd, be, de}

•••

At the end of the first scan,

if minsup = 80, remove ab, ad, ce, since count{ab, ad, ce} < 80</p>

Itemsets	Count		
{ab, ad, ce}	35		
{ <i>bd, be, de</i> }	298		
{ <i>yz, qs, wt</i> }	58		

Hash Table

<2> Direct Hashing and Pruning (DHP)

- DHP (Direct Hashing and Pruning): (J. Park, M. Chen, and P. Yu, SIGMOD'95)
- □ Hashing: Different itemsets may have the same hash value: v = hash(itemset)

 H_2

- □ 1st scan: When counting the 1-itemset, hash 2-itemset to calculate the bucket count
- Observation: A k-itemset cannot be frequent if its corresponding hashing bucket count is below the minsup threshold

Example:

	-							
Create hash table H_2	bucket address	0	1	2	3	4	5	6
using hash function	bucket count	2	2	4	2	2	4	4
$h(x, y) = ((order \ of \ x) \times 10$	bucket contents	{I1, I4}	{I1, I5}	{I2, I3}	{I2, I4}	{I2, I5}	{I1, I2}	{I1, I3}
+ (order of y)) mod 7		{I3, I5}	{I1, I5}	{I2, I3}	{I2, I4}	{I2, I5}	{I1, I2}	{I1, I3}
				{I2, I3}			{I1, I2}	{I1, I3}
				{I2, I3}			{I1, I2}	{I1, I3}

Figure 6.5: Hash table, H_2 , for candidate 2-itemsets: This hash table was generated by scanning the transactions of Table 6.1 while determining L_1 . If the minimum support count is, say, 3, then the itemsets in buckets 0, 1, 3, and 4 cannot be frequent and so they should not be included in C_2 .

<3> Exploring Vertical Data Format: ECLAT

- ECLAT (Equivalence Class Transformation): A depth-first search algorithm using set intersection [Zaki et al. @KDD'97]
- □ Tid-List: List of transaction-ids containing an itemset
- □ Vertical format: $t(e) = \{T_{10}, T_{20}, T_{30}\}; t(a) = \{T_{10}, T_{20}\}; t(ae) = \{T_{10}, T_{20}\}$
- Properties of Tid-Lists
 - t(X) = t(Y): X and Y always happen together (e.g., t(ac} = t(d))
 - □ $t(X) \subset t(Y)$: transaction having X always has Y (e.g., $t(ac) \subset t(ce)$)
- Deriving frequent patterns based on vertical intersections
- Using diffset to accelerate mining
 - Only keep track of differences of tids

□ t(e) = { T_{10} , T_{20} , T_{30} }, t(ce) = { T_{10} , T_{30} } → Diffset (ce, e) = { T_{20} }

A transaction DB in Horizontal Data Format

Tid	Itemset
10	a, c, d, e
20	a, b, e
30	b, c, e

The transaction DB in Vertical Data Format

Item	TidList
а	10, 20
b	20, 30
С	10, 30
d	10
е	10, 20, 30

<4> Mining Frequent Patterns by Pattern Growth

□ Apriori: A breadth-first search mining algorithm

- First find the complete set of frequent k-itemsets
- Then derive frequent (k+1)-itemset candidates
- Scan DB again to find true frequent (k+1)-itemsets

Two nontrivial costs:

- It may still need to generate a huge number of candidate sets. For example, if there are 10⁴ frequent 1-itemsets, the Apriori algorithm will need to generate more than 10⁷ candidate 2-itemsets.
- It may need to repeatedly scan the whole database and check a large set of candidates by pattern matching. It is costly to go over each transaction in the database to determine the support of the candidate itemsets.
<4> Mining Frequent Patterns by Pattern Growth

□ Apriori: A breadth-first search mining algorithm

- First find the complete set of frequent k-itemsets
- Then derive frequent (k+1)-itemset candidates
- Scan DB again to find true frequent (k+1)-itemsets
- Motivation for a different mining methodology
 - Can we mine the complete set of frequent patterns without such a costly generation process?
 - For a frequent itemset ρ, can subsequent search be confined to only those transactions that contain ρ?
 - A depth-first search mining algorithm?
- □ Such thinking leads to a frequent pattern (FP) growth approach:
 - **FPGrowth (J. Han, J. Pei, Y. Yin, "Mining Frequent Patterns without Candidate Generation," SIGMOD 2000)**

<4> High-level Idea of FP-growth Method

- Essence of frequent pattern growth (FPGrowth) methodology
 - Find frequent single items and partition the database based on each such single item pattern
 - Recursively grow frequent patterns by doing the above for each partitioned database (also called the pattern's conditional database)
 - To facilitate efficient processing, an efficient data structure, FP-tree, can be constructed
- Mining becomes
 - Recursively construct and mine (conditional) FP-trees
 - Until the resulting FP-tree is empty, or until it contains only one path—single path will generate all the combinations of its sub-paths, each of which is a frequent pattern

TID	Items in the Transaction	Ordered, frequent itemlist
100	{ <i>f</i> , <i>a</i> , <i>c</i> , <i>d</i> , <i>g</i> , <i>i</i> , <i>m</i> , <i>p</i> }	
200	$\{a, b, c, f, l, m, o\}$	
300	$\{b, f, h, j, o, w\}$	
400	$\{b, c, k, s, p\}$	
500	$\{a, f, c, e, l, p, m, n\}$	

1. Scan DB once, find single item frequent pattern:

Let min_support = 3

f:4, a:3, c:4, b:3, m:3, p:3

2. Sort frequent items in frequency descending order, f-list

F-list = f-c-a-b-m-p

TID	Items in the Transaction	Ordered, frequent itemlist
100	$\{f, a, c, d, g, i, m, p\}$	f, c, a, m, p
200	$\{a, b, c, f, l, m, o\}$	f, c, a, b, m
300	$\{b, f, h, j, o, w\}$	<i>f</i> , <i>b</i>
400	$\{b, c, k, s, p\}$	<i>c</i> , <i>b</i> , <i>p</i>
500	$\{a, f, c, e, l, p, m, n\}$	<i>f</i> , <i>c</i> , <i>a</i> , <i>m</i> , <i>p</i>

1. Scan DB once, find single item frequent pattern:

Let min_support = 3

f:4, a:3, c:4, b:3, m:3, p:3

2. Sort frequent items in frequency descending

order, f-list F-list = f-c-a-b-m-p

	TID	lte	ms in the Transaction	Order	dered, frequent itemlist		nlist	
	100		{f, a, c, d, g, i, m, p}		<i>f</i> , <i>c</i> ,	a, m, p		
	200		$\{a, b, c, f, l, m, o\}$		<i>f</i> , <i>c</i> ,	a, b, m		
	300		$\{b, f, h, j, o, w\}$		f	; b		Itemlist: "f c a m p"
	400		$\{b, c, k, s, p\}$		с,	<i>b</i> , <i>p</i>		
	500		{ <i>a</i> , <i>f</i> , <i>c</i> , <i>e</i> , <i>l</i> , <i>p</i> , <i>m</i> , <i>n</i> }		<i>f</i> , <i>c</i> ,	a, m, p		
1.	L. Scan DB once, find single item frequent pattern:			ŀ	leader Tak	ble		
Let min_supp		Let min_support = 3		Item	Frequency	header	f:1	
		f:4	, a:3, c:4, b:3, m:3, p:3					
2.	Sort	frequent	items in frequency desce	nding	f	4		c:I
	orde	r. f-list	F-list = f-c-a-b-m-n	0	С	4		
3	Scan	, DB again	construct FP-tree		a	3		a: 1
The frequent itemlist of each trans		ction is	b	3		$m \cdot 1$		
	inserted as a branch with shared sub-			m	3			
	br	anches m	nerged, counts accumulat	~ ed	р	3		$ \longrightarrow p:1 $

	TID	Items in the Transaction		Ordere	Ordered, frequent itemlist		nlist	
	100	100 $\{f, a, c, d, g, i, m, p\}$			<i>f</i> , <i>c</i> , <i>c</i>	a, m, p		
	200		{a, b, c, f, l, m, o}		<i>f</i> , <i>c</i> , <i>c</i>	a, b, m		
	300		$\{b, f, h, j, o, w\}$		ſ	с, b		After inserting the 2 nd frequent
	400		$\{b, c, k, s, p\}$		с,	<i>b</i> , <i>p</i>		itemlist "f, c, a, b, m"
	500		{ <i>a</i> , <i>f</i> , <i>c</i> , <i>e</i> , <i>l</i> , <i>p</i> , <i>m</i> , <i>n</i> }		<i>f</i> , <i>c</i> , <i>c</i>	a, m, p		
1.	L. Scan DB once, find single item frequent pattern: Header Table				ole			
	f:4, a:3, c:4, b:3, m:3, p:3			3	Item	Frequency	header	f:2
2	Sort	frequent	items in frequency de	escending	f	4		c:2
2.	orde	er. f-list	F-list - f-c-a-h-m-n		с	4		
ຊ	Scan	DB again	construct FP-tree		a	3		
5.	The frequent itemlist of each transaction is			nsaction is	b	3		$-\frac{1}{m\cdot l} - \frac{1}{b\cdot l}$
	inserted as a branch with shared sub-				m	3		
	br	anches m	erged, counts accum	ulated	р	3		$- \rightarrow p:1 \rightarrow m:1$

	TID	Items in the Transaction	Ordere	ed, fre	quent iter	nlist	
	100 $\{f, a, c, d, g, i, m, p\}$			<i>f</i> , <i>c</i> , <i>a</i> , <i>m</i> , <i>p</i>			
	200	$\{a, b, c, f, l, m, o\}$		<i>f</i> , <i>c</i> ,	a, b, m		
	300	$\{b, f, h, j, o, w\}$		j	f, b		After inserting all the
	400	$\{b, c, k, s, p\}$		с,	<i>b</i> , <i>p</i>		frequent itemlists
	500	$\{a, f, c, e, l, p, m, n\}$		<i>f</i> , <i>c</i> ,	a, m, p		
1.	Scan DB once, find single item frequent pattern: Header Table				ole		
		f:4, a:3, c:4, b:3, m:3, p:3		Item	Frequency	header	$f:4 \longrightarrow c:1$
2	Sort	frequent items in frequency descer	nding	f	4		$c:3 \rightarrow b:1 \rightarrow b:1$
۷.	orde	r. f-list $F_{-list} - f_{-c_{-}a_{-}b_{-}m_{-}n}$		с	4		
ว	Scan	DB again construct FP-tree		a	3		$ \Rightarrow a:3$
5.				b	3		

b:1

m:

m:2

p:2

3

3

m

р

The frequent itemlist of each transaction is inserted as a branch, with shared subbranches merged, counts accumulated

Mining FP-Tree: Divide and Conquer Based on Patterns and Data

- Pattern mining can be partitioned according to current patterns
 - Patterns containing p: p's conditional database: fcam:2, cb:1
 - p's conditional database (i.e., the database under the condition that p exists):
 - transformed prefix paths of item p
 - Patterns having m but no p: m's conditional database: fca:2, fcab:1



Mine Each Conditional Database Recursively

min_support = 3

Conditional Data Bases

item cond. data base

- c f:3
- a fc:3

81

- b fca:1, f:1, c:1
- m fca:2, fcab:1
- p fcam:2, cb:1



- For each conditional database
 - Mine single-item patterns
 - Construct its FP-tree & mine it

p's conditional DB: *fcam:2, cb:1 → c: 3*

m's conditional DB: *fca:2, fcab:1 → fca: 3*

b's conditional DB: *fca:1, f:1, c:1* $\rightarrow \phi$

Actually, for single branch FP-tree, all the frequent patterns can be generated in one shot

m: 3 fm: 3, cm: 3, am: 3 fcm: 3, fam:3, cam: 3 fcam: 3

A Special Case: Single Prefix Path in FP-tree

- Suppose a (conditional) FP-tree T has a shared single prefix-path P
- Mining can be decomposed into two parts
- Reduction of the single prefix path into one node
- Concatenation of the mining results of the two parts



 $\{\}$

 $a_1:n_1$

FPGrowth: Mining Frequent Patterns by Pattern Growth

- Essence of frequent pattern growth (FPGrowth) methodology
 - Find frequent single items and partition the database based on each such single item pattern
 - Recursively grow frequent patterns by doing the above for each partitioned database (also called the pattern's conditional database)
 - To facilitate efficient processing, an efficient data structure, FP-tree, can be constructed
- Mining becomes
 - Recursively construct and mine (conditional) FP-trees
 - Until the resulting FP-tree is empty, or until it contains only one path—single path will generate all the combinations of its sub-paths, each of which is a frequent pattern

Scaling FP-growth by Item-Based Data Projection

- What if FP-tree cannot fit in memory?—Do not construct FP-tree
 - "Project" the database based on frequent single items
 - Construct & mine FP-tree for each projected DB
- Parallel projection vs. partition projection
 - Parallel projection: Project the DB on each frequent item
 - Space costly, all partitions can be processed in parallel
 - Partition projection: Partition the DB in order
 - Passing the unprocessed parts to subsequent partitions





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

Basic Concepts

Efficient Pattern Mining Methods





Pattern Evaluation

Limitation of the Support-Confidence Framework

 \square Interestingness Measures: Lift and χ^2

Null-Invariant Measures

Comparison of Interestingness Measures

How to Judge if a Rule/Pattern Is Interesting?

- Pattern mining will generate a large set of patterns/rules
 - Not all the generated patterns/rules are interesting

How to Judge if a Rule/Pattern Is Interesting?

Pattern mining will generate a large set of patterns/rules

Not all the generated patterns/rules are interesting

□ Interestingness measures: Objective vs. Subjective

How to Judge if a Rule/Pattern Is Interesting?

- Pattern mining will generate a large set of patterns/rules
 - Not all the generated patterns/rules are interesting
- Interestingness measures: Objective vs. Subjective
 - Objective interestingness measures
 - Support, confidence, correlation, ...
 - Subjective interestingness measures:
 - Different users may judge interestingness differently
 - Let a user specify
 - Query-based: Relevant to a user's particular request
 - Judge against one's knowledge base
 - unexpected, freshness, timeliness

Limitation of the Support-Confidence Framework

□ Are *s* and *c* interesting in association rules: "A \Rightarrow B" [*s*, *c*]?

Limitation of the Support-Confidence Framework

- □ Are *s* and *c* interesting in association rules: "A \Rightarrow B" [*s*, *c*]?
- Example: Suppose one school may have the following statistics on # of students who may play basketball and/or eat cereal:

	play-basketball	not play-basketball	sum (row)
eat-cereal	400	350	750
not eat-cereal	200	50	250
sum(col.)	600	400	1000



Limitation of the Support-Confidence Framework

- □ Are *s* and *c* interesting in association rules: "A \Rightarrow B" [*s*, *c*]?
- Example: Suppose one school may have the following statistics on # of students who may play basketball and/or eat cereal:

	play-basketball	not play-basketball	sum (row)
eat-cereal	400	350	750
not eat-cereal	200	50	250
sum(col.)	600	400	1000



□ Association rule mining may generate the following:

■ *play-basketball* \Rightarrow *eat-cereal* [40%, 66.7%] (higher s & c)

- But this strong association rule is misleading: The overall % of students eating cereal is 75% > 66.7%, a more telling rule:
 - \neg play-basketball \Rightarrow eat-cereal [35%, 87.5%] (high s & c)

Interestingness Measure: Lift

Measure of dependent/correlated events: lift

$$lift(B,C) = \frac{c(B \to C)}{s(C)} = \frac{P(C|B)}{P(C)} = \frac{P(B \cup C)}{P(B)P(C)}$$

Lift is more telling than s & c

	В	¬Β	Σ _{row}
С	400	350	750
¬C	200	50	250
Σ _{col} .	600	400	1000

Interestingness Measure: Lift

Measure of dependent/correlated events: lift

$$lift(B,C) = \frac{c(B \rightarrow C)}{s(C)} = \frac{P(C|B)}{P(C)} = \frac{P(B \cup C)}{P(B)P(C)}$$

- Lift(B, C) may tell how B and C are correlated
 - \Box Lift(B, C) = 1: B and C are independent
 - \square > 1: positively correlated
 - \Box < 1: negatively correlated

Lift is more telling than s & c

	В	¬Β	Σ _{row}
С	400	350	750
٦C	200	50	250
Σ _{col.}	600	400	1000

Interestingness Measure: Lift

Measure of dependent/correlated events: lift

$$lift(B,C) = \frac{c(B \rightarrow C)}{s(C)} = \frac{P(C|B)}{P(C)} = \frac{P(B \cup C)}{P(B)P(C)}$$

- □ Lift(B, C) may tell how B and C are correlated
 - \Box Lift(B, C) = 1: B and C are independent
 - \square > 1: positively correlated
 - □ < 1: negatively correlated

In our example,

$$lift(B,C) = \frac{400/1000}{600/1000 \times 750/1000} = 0.89$$

$$lift(B,\neg C) = \frac{200/1000}{600/1000 \times 250/1000} = 1.33$$

□ Thus, B and C are negatively correlated since lift(B, C) < 1;

□ B and \neg C are positively correlated since lift(B, \neg C) > 1

Lift is more telling than s & c

	В	¬Β	Σ _{row}
С	400	350	750
٦C	200	50	250
Σ _{col} .	600	400	1000

Interestingness Measure: χ^2

 \Box Another measure to test correlated events: χ^2

$$\chi^{2} = \sum \frac{(Observed - Expected)^{2}}{Expected}$$

		В	¬B	Σ _{row}
С	400 (450)		350 (300)	750
¬C	200	(150)	50 (100)	250
Σ_{col}		600	400	1000

Expected value

Observed value

Interestingness Measure: χ^2

 \square Another measure to test correlated events: χ^2

$$\chi^{2} = \sum \frac{(Observed - Expected)^{2}}{Expected}$$

□ For the table on the right,

$$\chi^{2} = \frac{(400 - 450)^{2}}{450} + \frac{(350 - 300)^{2}}{300} + \frac{(200 - 150)^{2}}{150} + \frac{(50 - 100)^{2}}{100} = 55.56$$

Observed value

Expected value

¬Β

350 (300)

50 (100)

400

 Σ_{row}

750

250

1000

Β

400 (450)

20((150)

600

С

¬С

 Σ_{col}

- □ By consulting a table of critical values of the χ^2 distribution, one can conclude that the chance for B and C to be independent is very low (< 0.01)
- χ²-test shows B and C are negatively correlated since the expected value is
 450 but the observed is only 400
- $\hfill \hfill \hfill$

Lift and χ^2 : Are They Always Good Measures?

- Null transactions: Transactions that contain neither B nor C
- Let's examine the new dataset D
 - BC (100) is much rarer than B¬C (1000) and ¬BC (1000), but there are many ¬B¬C (100000)
 - Unlikely B & C will happen together!
- But, Lift(B, C) = 8.44 >> 1 (Lift shows B and C are strongly positively correlated!)
- $\square \chi^2 = 670$: Observed(BC) >> expected value (11.85)
- Too many null transactions may "spoil the soup"!

		В	¬B	Σ_{row}			
	С	100	1000	1100			
	٦C	1000	100000	101000			
	$\Sigma_{col.}$	1100 🧳	101000	102100			
•							

Contingency table with expected values added

null transactions

	В	¬В	Σ_{row}	
С	100 (11.85)	1000	1100	
٦C	1000 (988.15)	100000	101000	
Σ _{col.}	1100	101000	102100	



Interestingness Measures & Null-Invariance

- Null invariance: Value does not change with the # of null-transactions
- □ A few interestingness measures: Some are null invariant

Measure	Definition	Range	Null-Invariant?]
$\chi^2(A,B)$	$\sum_{i,j} \frac{(e(a_i,b_j)-o(a_i,b_j))^2}{e(a_i,b_j)}$	$[0,\infty]$	No	X ² and lift are not
Lift(A, B)	$\frac{s(A \cup B)}{s(A) \times s(B)}$	$[0,\infty]$	No	
Allconf(A, B)	$\frac{s(A \cup B)}{max\{s(A), s(B)\}}$	[0,1]	Yes] /
Jaccard(A, B)	$\frac{s(A \cup B)}{s(A) + s(B) - s(A \cup B)}$	[0, 1]	Yes	Jaccard, consine,
Cosine(A, B)	$\frac{s(A \cup B)}{\sqrt{s(A) \times s(B)}}$	[0, 1]	Yes	and Kulczynski are
Kulczynski(A, B)	$\frac{1}{2}\left(\frac{s(A\cup B)}{s(A)} + \frac{s(A\cup B)}{s(B)}\right)$	[0, 1]	Yes	null-invariant measures
MaxConf(A, B)	$\max\{\frac{s(A\cup B)}{s(A)}, \frac{s(A\cup B)}{s(B)}\}$	[0, 1]	Yes	

Null Invariance: An Important Property

- Why is null invariance crucial for the analysis of massive transaction data?
 - Many transactions may contain neither milk nor coffee!

	milk	$\neg milk$	Σ_{row}
coffee	mc	$\neg mc$	c
$\neg coffee$	$m \neg c$	$\neg m \neg c$	$\neg c$
Σ_{col}	\overline{m}	$\neg m$	Σ

milk vs. coffee contingency table

- □ Lift and χ^2 are not null-invariant: not good to evaluate data that contain too many or too few null transactions!
- Many measures are not null-invariant!

Null-transactions

Σ_{col} m	$\neg m$	Σ	w.r.t. m	and c		
Data set	mc	$\neg mc$	$m \neg c$	$n\neg c$	χ^2	Lift
D_1	10,000	1,000	1,000	100,000	90557	9.26
D_2	10,000	1,000	1,000	100	0	1
D_3	100	1,000	1,000	100,000	670	8.44
D_4	1,000	1,000	1,000	100,000	24740	25.75
D_5	1,000	100	10,000	100,000	8173	9.18
D_6	1,000	10	100,000	100,000	965	1.97

Comparison of Null-Invariant Measures

- Not all null-invariant measures are created equal
- Which one is better?
 - **D**₄ $-D_6$ differentiate the null-invariant measures
 - Kulc (Kulczynski 1927) holds firm and is in balance of both directional implications

2-variable contingency table

	milk	$\neg milk$	Σ_{row}
$co\!f\!fee$	mc	$\neg mc$	c
$\neg coffee$	$m \neg c$	$\neg m \neg c$	$\neg c$
Σ_{col}	m	$\neg m$	Σ

				All 5 are nu	ll-invariant		· · ·		•
Data set	mc	$\neg mc$	$m \neg c$	$\neg m \neg c$	AllConf	Jaccard	Cosine	Kulc	MaxConf
D_1	10,000	1,000	1,000	100,000	0.91	0.83	0.91	0.91	0.91
D_2	10,000	1,000	1,000	100	0.91	0.83	0.91	0.91	0.91
D_3	100	1,000	1,000	100,000	0.09	0.05	0.09	0.09	0.09
D_4	1,000	1,000	1,000	100,000	0.5	0.33	0.5	0.5	0.5
D_5	1,000	100	10,000	100,000	0.09	0.09	0.29	0.5	0.91
D_6	1,000	10	100,000	<u>100,000</u>	0.01	0.01	0.10	0.5	0.99

Subtle: They disagree on those cases

Imbalance Ratio with Kulczynski Measure

□ IR (Imbalance Ratio): measure the imbalance of two itemsets A and B in rule implications:

$$IR(A,B) = \frac{|s(A) - s(B)|}{s(A) + s(B) - s(A \cup B)}$$

- Kulczynski and Imbalance Ratio (IR) together present a clear picture for all the three datasets D₄ through D₆
 - \square D₄ is neutral & balanced; D₅ is neutral but imbalanced
 - \square D₆ is neutral but very imbalanced

Data set	mc	$\neg mc$	$m \neg c$	$\neg m \neg c$	Jaccard	Cosine	Kulc	IR
D_1	10,000	1,000	1,000	100,000	0.83	0.91	0.91	0
D_2	10,000	1,000	1,000	100	0.83	0.91	0.91	0
D_3	100	1,000	1,000	100,000	0.05	0.09	0.09	0
D_4	1,000	1,000	1,000	100,000	0.33	0.5	≤ 0.5	0
D_5	1,000	100	10,000	100,000	0.09	0.29	$\bigcirc 0.5$	0.89
D_6	1,000	10	100,000	100,000	0.01	0.10	$\bigcirc 0.5$	0.99

What Measures to Choose for Effective Pattern Evaluation?

- Null value cases are predominant in many large datasets
 - Neither milk nor coffee is in most of the baskets; neither Mike nor Jim is an author in most of the papers;
- □ *Null-invariance* is an important property
- Lift, χ² and cosine are good measures if null transactions are not predominant
 Otherwise, *Kulczynski* + *Imbalance Ratio* should be used to judge the interestingness of a pattern

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

Basic Concepts

Efficient Pattern Mining Methods

Pattern Evaluation

