



06 pattern Analysis

Dr. N. Abdolvand

Data Mining

Data Mining: Concepts and Techniques

(3rd ed.)

— Chapter 6 —

Jiawei Han, Micheline Kamber, and Jian Pei
University of Illinois at Urbana-Champaign &
Simon Fraser University

©2013 Han, Kamber & Pei. All rights reserved.

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

- ◆ Basic Concepts
- ◆ Frequent Itemset Mining Methods
- ◆ Which Patterns Are Interesting?—Pattern Evaluation Methods
- ◆ Summary

3

Case Study



- ◆ Roger is a city manager for a medium-sized, but steadily growing, city. The city has limited resources, and like most municipalities, there are more needs than there are resources. He feels like the citizens in the community are fairly active in various community organizations, and believes that he may be able to get a number of groups to work together to meet some of the needs in the community. He knows there are churches, social clubs, hobby enthusiasts and other types of groups in the community. What he doesn't know is if there are connections between the groups that might enable natural collaborations between two or more groups that could work together on projects around town. He decides that before he can begin asking community organizations to begin working together and to accept responsibility for projects, he needs to find out if there are any existing associations between the different types of groups in the area.

Dr. N. Abdolvand

ORGANIZATIONAL UNDERSTANDING

- ◆ Roger's goal is to identify and then try to take advantage of existing connections in his local community to get some work done that will benefit the entire community. He knows of many of the organizations in town, has contact information for them and is even involved in some of them himself. His family is involved in an even broader group of organizations, so he understands on a personal level the diversity of groups and their interests. Because people he and his family knows are involved in other groups around town, he is aware in a more general sense of many different types of organizations, their interests, objectives and potential contributions. He knows that to start, his main concern is finding types of organizations that seem to be connected with one another. Identifying individuals to work with at each church, social club or political organization will be overwhelming without first categorizing the organizations into groups and looking for associations between the groups. Only once he's checked for existing connections will he feel ready to begin contacting people and asking them to use their cross-organizational contacts and take on project ownership. His first need is to find where such associations exist



Dr. N. Abdolvand

attributes

- ◆ **Elapsed-Time:** This is the amount of time each respondent spent completing our survey. It is expressed in decimal minutes
- ◆ **Time-in-Community:** This question on the survey asked the person if they have lived in the area for 0-2 years, 3-9 years, or 10+ years; and is recorded in the data set as Short, Medium, or Long respectively.
- ◆ **Gender:** The survey respondent's gender.
- ◆ **Working:** A yes/no column
- ◆ **Age:** The survey respondent's age in years.
- ◆ **Family:** A yes/no column indicating whether or not the respondent is currently a member of a family-oriented community organization, such as Big Brothers/Big Sisters, childrens' recreation or sports leagues, genealogy groups, etc.
- ◆ **Hobbies:** A yes/no column indicating whether or not the respondent is currently a member of a hobby-oriented community organization, such as amateur radio, outdoor recreation, motorcycle or bicycle riding, etc.



Dr. N. Abdolvand

attributes

- ◆ **Social-Club:** A yes/no column indicating whether or not the respondent is currently a member of a community social organization, such as Rotary International, Lion's Club, etc.
- ◆ **Political:** A yes/no column indicating whether or not the respondent is currently a member of a political organization with regular meetings in the community, such as a political party, a grass-roots action group, a lobbying effort, etc.
- ◆ **Professional:** A yes/no column indicating whether or not the respondent is currently a member of a professional organization with local chapter meetings, such as a chapter of a law or medical society, a small business owner's group, etc.
- ◆ **Religious:** A yes/no column indicating whether or not the respondent is currently a member of a church in the community.
- ◆ **Support-Group:** A yes/no column indicating whether or not the respondent is currently a member of a support-oriented community organization, such as Alcoholics Anonymous, an anger management group, etc.



Dr. N. Abdolvand

Supervised Versus Unsupervised Methods

- ◆ **Supervised learning:** is typically done in the context of classification, when we want to map input to output labels, or regression, when we want to map input to a continuous output.
- ◆ **unsupervised learning:** here we wish to learn the inherent structure of our data without using explicitly-provided labels.

	<i>Supervised Learning</i>	<i>Unsupervised Learning</i>
<i>Discrete</i>	classification or categorization	clustering
<i>Continuous</i>	regression	dimensionality reduction



Dr. N. Abdolvand

What Is Frequent Pattern Analysis?

- ◆ Frequent pattern: a pattern (a set of items, subsequences, substructures, etc.) that occurs frequently in a data set
- ◆ First proposed by Agrawal, Imielinski, and Swami [AIS93] in the context of frequent itemsets and association rule mining
- ◆ Motivation: Finding inherent regularities in data
 - What products were often purchased together?— Beer and diapers?!
 - What are the subsequent purchases after buying a PC?
 - What kinds of DNA are sensitive to this new drug?
 - Can we automatically classify web documents?
- ◆ Applications
 - Basket data analysis, cross-marketing, catalog design, sale campaign analysis, Web log (click stream) analysis, and DNA sequence analysis.

9

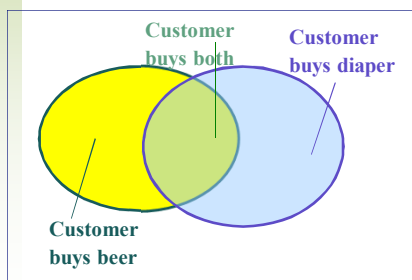
Why Is Freq. Pattern Mining Important?

- ◆ Freq. pattern: An intrinsic and important property of datasets
- ◆ Foundation for many essential data mining tasks
 - Association, correlation, and causality analysis
 - Sequential, structural (e.g., sub-graph) patterns
 - Pattern analysis in spatiotemporal, multimedia, time-series, and stream data
 - Classification: discriminative, frequent pattern analysis
 - Cluster analysis: frequent pattern-based clustering
 - Data warehousing: iceberg cube and cube-gradient
 - Semantic data compression: fascicles
 - Broad applications

10

Basic Concepts: Frequent Patterns

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

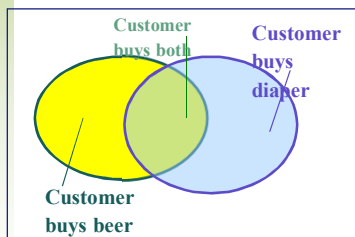


- ◆ **itemset**: A set of one or more items
- ◆ **k-itemset** $X = \{x_1, \dots, x_k\}$
- ◆ **(absolute) support**, or, **support count** of X : Frequency or occurrence of an itemset X
- ◆ **(relative) support**, s , is the fraction of transactions that contains X (i.e., the probability that a transaction contains X)
- ◆ An itemset X is **frequent** if X 's support is no less than a **minsup** threshold

11

Basic Concepts: Association 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



- ◆ Find all the rules $X \rightarrow Y$ with minimum support and confidence
 - **support**, s , probability that a transaction contains $X \cup Y$
 - **confidence**, c , conditional probability that a transaction having X also contains Y
- Let $\text{minsup} = 50\%$, $\text{minconf} = 50\%$
 Freq. Pat.: Beer:3, Nuts:3, Diaper:4, Eggs:3, {Beer, Diaper}:3
- Association rules: (many more!)
 - $\text{Beer} \rightarrow \text{Diaper}$ (60%, 100%)
 - $\text{Diaper} \rightarrow \text{Beer}$ (60%, 75%)

12

Associations Among Facebook Likes

246

- ◆ Although finding associations is often used with market basket data—and sometimes is even called **market-basket analysis**—the technique is much more general. We have data on the things that were “Liked” by a large collection of users of Facebook. By analogy to market basket data, we can consider each of these users to have a “basket” of Likes, by aggregating all the Likes of each user.

Like	Like
<i>Lord Of The Rings</i>	Wikileaks
One Manga	Beethoven
Science	NPR
Psychology	<i>Spirited Away</i>
<i>The Big Bang Theory</i>	Running
Paulo Coelho	Roger Federer
<i>The Daily Show</i>	<i>Star Trek</i>
<i>Lost</i>	Philosophy
<i>Lie to Me</i>	<i>The Onion</i>
<i>How I Met Your Mother</i>	<i>The Colbert Report</i>
<i>Doctor Who</i>	<i>Star Trek</i>
<i>Howl's Moving Castle</i>	Sheldon Cooper
<i>Tron</i>	<i>Fight Club</i>
Angry Birds	<i>Inception</i>
<i>The Godfather</i>	<i>Weeds</i>

Dr. N. Abdolvand

- ◆ What we would like to do is to look to see what are the Likes that give strong evidence lifts for “high IQ,” or more specifically for scoring high on an IQ test. Taking a sample of the Facebook population, if we define our target variable as the binary variable $IQ > 130$, about 14% of the sample is positive (has $IQ > 130$).
- ◆ “Liking” *Fight Club*, *Star Trek*, or *Sheldon Cooper* on Facebook each increases by about 30% the estimation of the probability that you have a high IQ. If you were to Like all three of those, it would more than double the estimate that you have a high IQ.
- ◆ It means clicking “Like” on Facebook’s page called **The Lord of the Rings** is a strong indication that I have very high IQ



Dr. N. Abdolvand

Scalable Frequent Itemset Mining Methods

- ◆ Apriori: A Candidate Generation-and-Test Approach
- ◆ Improving the Efficiency of Apriori
- ◆ FPGrowth: A Frequent Pattern-Growth Approach
- ◆ ECLAT: Frequent Pattern Mining with Vertical Data Format

17

The Downward Closure Property and Scalable Mining Methods

- ◆ The downward closure property of frequent patterns
 - Any subset of a frequent itemset must be frequent
 - If {beer, diaper, nuts} is frequent, so is {beer, diaper}
 - i.e., every transaction having {beer, diaper, nuts} also contains {beer, diaper}
- ◆ Scalable mining methods: Three major approaches
 - Apriori (Agrawal & Srikant@VLDB'94)
 - Freq. pattern growth (FPgrowth—Han, Pei & Yin @SIGMOD'00)
 - Vertical data format approach (Charm—Zaki & Hsiao @SDM'02)

18

Apriori: A Candidate Generation & Test Approach

- ◆ Apriori pruning principle: If there is any itemset which is infrequent, its superset should not be generated/tested! (Agrawal & Srikant @VLDB'94, Mannila, et al. @ KDD' 94)
- ◆ Method:
 - Initially, scan DB once to get frequent 1-itemset
 - Generate length (k+1) candidate itemsets from length k frequent itemsets
 - Test the candidates against DB
 - Terminate when no frequent or candidate set can be generated

19

The Apriori Algorithm—An Example

Database TDB

Tid	Items
10	A, C, D
20	B, C, E
30	A, B, C, E
40	B, E

$Sup_{min} = 2$

1st scan

Itemset	sup
{A}	2
{B}	3
{C}	3
{D}	1
{E}	3

Itemset	sup
{A}	2
{B}	3
{C}	3
{E}	3

2nd scan

Itemset	sup
{A, B}	1
{A, C}	2
{A, E}	1
{B, C}	2
{B, E}	3
{C, E}	2

Itemset	sup
{A, C}	2
{B, C}	2
{B, E}	3
{C, E}	2

3rd scan

Itemset	sup
{B, C, E}	2

Itemset	sup
{B, C, E}	2

20

Further Improvement of the Apriori Method

- ◆ Major computational challenges
 - Multiple scans of transaction database
 - Huge number of candidates
 - Tedious workload of support counting for candidates
- ◆ Improving Apriori: general ideas
 - Reduce passes of transaction database scans
 - Shrink number of candidates
 - Facilitate support counting of candidates

21

Advantages of the Pattern Growth Approach

- ◆ Divide-and-conquer:
 - Decompose both the mining task and DB according to the frequent patterns obtained so far
 - Lead to focused search of smaller databases
- ◆ Other factors
 - No candidate generation, no candidate test
 - Compressed database: FP-tree structure
 - No repeated scan of entire database
 - Basic ops: counting local freq items and building sub FP-tree, no pattern search and matching
- ◆ A good open-source implementation and refinement of FPGrowth
 - FPGrowth+ (Grahne and J. Zhu, FIMI'03)

24

Mining Frequent Itemsets without Candidate Generation

- ◆ frequent-pattern (FP) growth
 - Devide & Conquer strategy
 - Compresses the database representing frequent items into a frequent-pattern tree
 - then divides the compressed database into a set of *conditional databases each* associated with one frequent item or “pattern fragment,” and mines each such database separately.
 -

Dr. N. Abdolvand

Table 5.1 Transactional data for an *AllElectronics* branch.

TID	List of item IDs
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

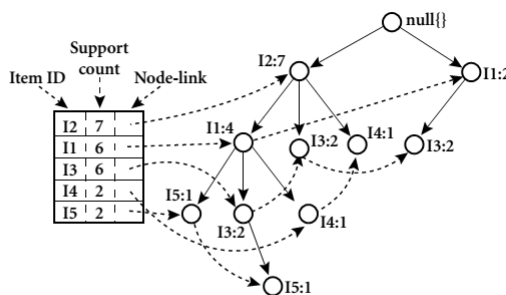


Figure 5.7 An FP-tree registers compressed, frequent pattern information.

Table 5.2 Mining the FP-tree by creating conditional (sub-)pattern bases.

Item	Conditional Pattern Base	Conditional FP-tree	Frequent Patterns Generated
I5	{{I2, I1: 1}, {I2, I1, I3: 1}}	$\langle I2: 2, I1: 2 \rangle$	{I2, I5: 2}, {I1, I5: 2}, {I2, I1, I5: 2}
I4	{{I2, I1: 1}, {I2: 1}}	$\langle I2: 2 \rangle$	{I2, I4: 2}
I3	{{I2, I1: 2}, {I2: 2}, {I1: 2}}	$\langle I2: 4, I1: 2 \rangle, \langle I1: 2 \rangle$	{I2, I3: 4}, {I1, I3: 4}, {I2, I1, I3: 2}
I1	{{I2: 4}}	$\langle I2: 4 \rangle$	{I2, I1: 4}

Dr. N. Abdolvand

Mining Frequent Itemsets Using Vertical Data Format

◆ ECLAT

Dr. N. Abdolvand

Table 5.1 Transactional data for an *AlI*Electronics branch.

<i>TID</i>	<i>List of item IDs</i>
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

Table 5.3 The vertical data format of the transaction data set *D* of Table 5.1.

<i>itemset</i>	<i>TID_set</i>
I1	{T100, T400, T500, T700, T800, T900}
I2	{T100, T200, T300, T400, T600, T800, T900}
I3	{T300, T500, T600, T700, T800, T900}
I4	{T200, T400}
I5	{T100, T800}

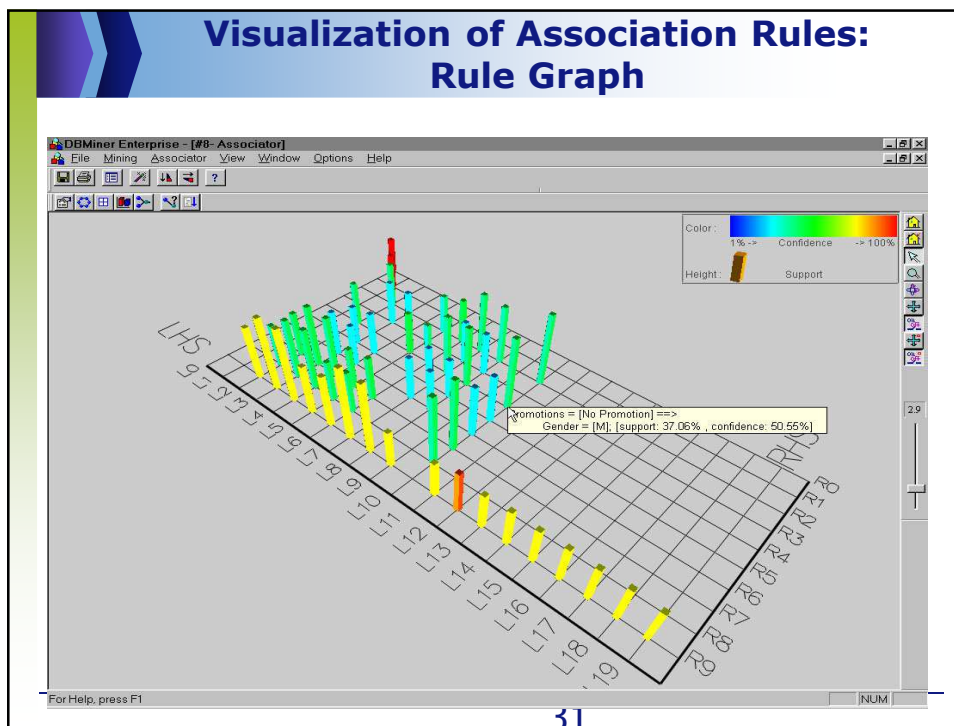
Table 5.4 The 2-itemsets in vertical data format.

<i>itemset</i>	<i>TID_set</i>
{I1, I2}	{T100, T400, T800, T900}
{I1, I3}	{T500, T700, T800, T900}
{I1, I4}	{T400}
{I1, I5}	{T100, T800}
{I2, I3}	{T300, T600, T800, T900}
{I2, I4}	{T200, T400}
{I2, I5}	{T100, T800}
{I3, I5}	{T800}

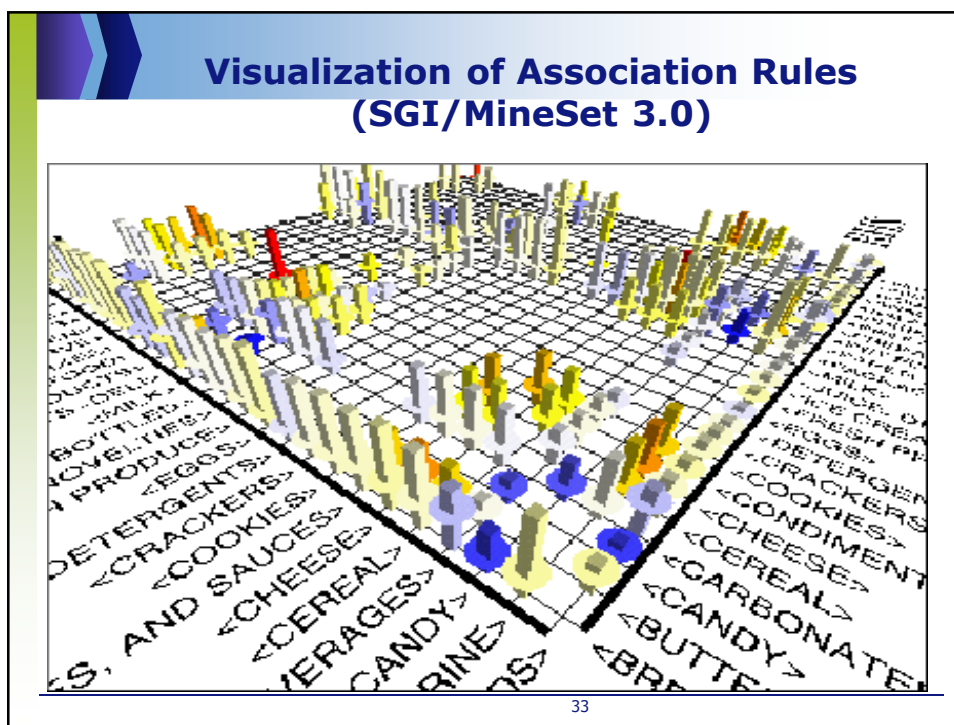
Table 5.5 The 3-itemsets in vertical data format.

<i>itemset</i>	<i>TID_set</i>
{I1, I2, I3}	{T800, T900}
{I1, I2, I5}	{T100, T800}

Dr. N. Abdolvand



31



33

Computational Complexity of Frequent Itemset Mining

- ◆ How many itemsets are potentially to be generated in the worst case?
 - The number of frequent itemsets to be generated is sensitive to the minsup threshold
 - When minsup is low, there exist potentially an exponential number of frequent itemsets
 - The worst case: MN where M : # distinct items, and N : max length of transactions
- ◆ The worst case complexity vs. the expected probability
 - Ex. Suppose Walmart has 104 kinds of products
 - The chance to pick up one product 10^{-4}
 - The chance to pick up a particular set of 10 products: $\sim 10^{-40}$
 - What is the chance this particular set of 10 products to be frequent 103 times in 109 transactions?

34

Interestingness Measure: Correlations (Lift)

- ◆ $play\ basketball \Rightarrow eat\ cereal$ [40%, 66.7%] is misleading
 - The overall % of students eating cereal is 75% > 66.7%.
- ◆ $play\ basketball \Rightarrow not\ eat\ cereal$ [20%, 33.3%] is more accurate, although with lower support and confidence
- ◆ Measure of dependent/correlated events: lift

$$lift = \frac{P(A \cup B)}{P(A)P(B)}$$

$$lift(B, C) = \frac{2000 / 5000}{3000 / 5000 * 3750 / 5000} = 0.89$$

$$lift(B, \neg C) = \frac{1000 / 5000}{3000 / 5000 * 1250 / 5000} = 1.33$$

	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum (col.)	3000	2000	5000

35

Strong Rules Are Not Necessarily Interesting: An Example

$$\text{buys}(X, \text{"computer games"}) \Rightarrow \text{buys}(X, \text{"videos"}) \quad [\text{support} = 40\%, \text{confidence} = 66\%]$$

Rule (5.21) is a strong association rule and would therefore be reported, since its support value of $\frac{4,000}{10,000} = 40\%$ and confidence value of $\frac{4,000}{6,000} = 66\%$ satisfy the minimum support and minimum confidence thresholds, respectively. However, Rule (5.21) is misleading because the probability of purchasing videos is 75%, which is even larger than 66%. In fact, computer games and videos are negatively associated because the purchase of one of these items actually decreases the likelihood of purchasing the other. Without fully understanding this phenomenon, we could easily make unwise business decisions based on Rule (5.21). ■

can see that the probability of purchasing a computer game is $P(\{\text{game}\}) = 0.60$, the probability of purchasing a video is $P(\{\text{video}\}) = 0.75$, and the probability of purchasing both is $P(\{\text{game}, \text{video}\}) = 0.40$. By Equation (5.23), the lift of Rule (5.21) is $P(\{\text{game}, \text{video}\}) / (P(\{\text{game}\}) \times P(\{\text{video}\})) = 0.40 / (0.60 \times 0.75) = 0.89$. Because this value is less than 1, there is a negative correlation between the occurrence of $\{\text{game}\}$ and $\{\text{video}\}$. The numerator is the likelihood of a customer purchasing both, while the

Dr. N. Abdolvand

Are lift and χ^2 Good Measures of Correlation?

- “Buy walnuts \Rightarrow buy milk [1%, 80%]” is misleading if 85% of customers buy milk
- Support and confidence are not good to indicate correlations
- Over 20 interestingness measures have been proposed (see Tan, Kumar, Sritastava @KDD'02)
- Which are good ones?

symbol	measure	range	formula
ϕ	ϕ -coefficient	-1 ... 1	$\frac{P(A,B) - P(A)P(B)}{\sqrt{(P(A)P(B)(1-P(A))(1-P(B)))}}$
Q	Yule's Q	-1 ... 1	$\frac{P(A,B)P(\bar{A},\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A},\bar{B}) + P(A,\bar{B})P(\bar{A},B)}$
Y	Yule's Y	-1 ... 1	$\frac{\sqrt{P(A,B)P(\bar{A},\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A},\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}}$
k	Cohen's k	-1 ... 1	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
PS	Piatetsky-Shapiro's	-0.25 ... 0.25	$P(A,B) - P(A)P(B)$
F	Certainty factor	-1 ... 1	$\max(\frac{P(A,B) - P(A)P(B)}{1 - P(B)}, \frac{P(A,B) - P(A)P(B)}{1 - P(A)})$
AV	added value	-0.5 ... 1	$\max(P(B A) - P(B), P(A B) - P(A))$
K	Klosgen's Q	0.33 ... 0.38	$\sqrt{\frac{P(A,B) \max(P(B A) - P(B), P(A B) - P(A))}{\sum_i \max_k P(A_i, B_i) + \sum_i \max_k P(A_i, \bar{B}_i) - \max_k P(A_i) - \max_k P(\bar{B}_i)}}$
g	Goodman-kruskal's	0 ... 1	$\frac{\sum_i \sum_j P(A_i, B_j) - \max_k P(A_i)P(B_j)}{1 - \max_k P(A_i) - \max_k P(B_j)}$
M	Mutual Information	0 ... 1	$\frac{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)} + \sum_i \sum_j P(A_i, \bar{B}_j) \log \frac{P(A_i, \bar{B}_j)}{P(A_i)P(\bar{B}_j)}}$
J	J-Measure	0 ... 1	$\max(P(A, B) \log \frac{P(A, B)}{P(A)P(B)} + P(A, \bar{B}) \log \frac{P(A, \bar{B})}{P(A)P(\bar{B})}, P(A, B) \log \frac{P(A, B)}{P(A)P(B)} + P(\bar{A}, B) \log \frac{P(\bar{A}, B)}{P(\bar{A})P(B)})$
G	Gini index	0 ... 1	$\max(P(A)[P(B A)^2 + P(\bar{B} A)^2] + P(\bar{A})[P(B \bar{A})^2 + P(\bar{B} \bar{A})^2] - P(B)^2 - P(\bar{B})^2, P(B)[P(A B)^2 + P(\bar{A} B)^2] + P(\bar{B})[P(A \bar{B})^2 + P(\bar{A} \bar{B})^2] - P(A)^2 - P(\bar{A})^2)$
s	support	0 ... 1	$\frac{P(A, B)}{N}$
c	confidence	0 ... 1	$\frac{\max(P(B A), P(A B))}{P(A)}$
L	Laplace	0 ... 1	$\max(\frac{P(A, B) + 1}{N(A) + 2}, \frac{P(A, B) + 1}{N(A) + 2})$
IS	Cosine	0 ... 1	$\frac{P(A, B)}{\sqrt{P(A)P(B)}}$
γ	coherence (Jaccard)	0 ... 1	$\frac{P(A, B)}{P(A) + P(B) - P(A, B)}$
α	all_confidence	0 ... 1	$\frac{P(A, B)}{\max(P(A), P(B))}$
o	odds ratio	0 ... ∞	$\frac{P(A, B)P(\bar{A}, \bar{B})}{P(A, \bar{B})P(\bar{A}, B)}$
V	Conviction	0.5 ... ∞	$\max(\frac{P(A, B)}{P(A, \bar{B})}, \frac{P(B, \bar{A})}{P(\bar{A}, B)})$
λ	lift	0 ... ∞	$\frac{P(A, B)}{P(A)P(B)}$
S	Collective strength	0 ... ∞	$\frac{P(A, B) + P(\bar{A}, \bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A, B) - P(\bar{A}, \bar{B})}$
χ^2	χ^2	0 ... ∞	$\frac{\sum_i (P(A_i, B_i) - E_i)^2}{E_i}$

Null-Invariant Measures

Table 6: Properties of interestingness measures. Note that none of the measures satisfies all the properties.

Symbol	Measure	Range	P1	P2	P3	O1	O2	O3	O3'	O4
ϕ	ϕ -coefficient	-1...0...1	Yes	Yes	Yes	Yes	No	Yes	Yes	No
λ	Goodman-Kruskal's	0...1	Yes	No	No	Yes	No	No*	Yes	No
α	odds ratio	0...1... ∞	Yes*	Yes	Yes	Yes	Yes	Yes*	Yes	No
Q	Yule's Q	-1...0...1	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No
Y	Yule's Y	-1...0...1	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No
κ	Cohen's	-1...0...1	Yes	Yes	Yes	Yes	No	No	Yes	No
M	Mutual Information	0...1	Yes	Yes	Yes	No**	No	No*	Yes	No
J	J-Measure	0...1	Yes	No	No	No**	No	No	No	No
G	Gini index	0...1	Yes	No	No	No**	No	No*	Yes	No
s	Support	0...1	No	Yes	No	Yes	No	No	No	No
c	Confidence	0...1	No	Yes	No	No**	No	No	No	Yes
L	Laplace	0...1	No	Yes	No	No**	No	No	No	No
V	Conviction	0.5...1... ∞	No	Yes	No	No**	No	No	Yes	No
I	Interest	0...1... ∞	Yes*	Yes	Yes	Yes	No	No	No	No
IS	Cosine	$0 \dots \sqrt{P(A, B)}$	No	Yes	Yes	Yes	No	No	No	Yes
PS	Piatetsky-Shapiro's	-0.25...0...0.25	Yes	Yes	Yes	Yes	No	Yes	No	No
F	Certainty factor	-1...0...1	Yes	Yes	Yes	No**	No	No	Yes	No
AV	Added value	-0.5...0...1	Yes	Yes	Yes	No**	No	No	No	No
S	Collective strength	0...1... ∞	No	Yes	Yes	Yes	No	Yes*	Yes	No
ζ	Jaccard	0...1	No	Yes	Yes	Yes	No	No	No	Yes
K	Klorgen's	$(\frac{2}{\sqrt{3}} - 1)^{1/2} [2 - \sqrt{3} - \frac{1}{\sqrt{3}}] \dots 0 \dots \frac{2}{3\sqrt{3}}$	Yes	Yes	Yes	No**	No	No	No	No

where: P1: $O(M) = 0$ if $det(M) = 0$, i.e., whenever A and B are statistically independent.
 P2: $O(M_2) > O(M_1)$ if $M_2 = M_1 + [k \ -k; \ -k \ k]$.
 P3: $O(M_2) < O(M_1)$ if $M_2 = M_1 + [0 \ k; \ 0 \ -k]$ or $M_2 = M_1 + [0 \ 0; \ k \ -k]$.
 O1: Property 1: Symmetry under variable permutation.
 O2: Property 2: Row and Column scaling invariance.
 O3: Property 3: Antisymmetry under row or column permutation.
 O3': Property 4: Inversion invariance.
 O4: Property 5: Null invariance.
 Yes*: Yes if measure is normalized.
 No*: Symmetry under row or column permutation.
 No**: No unless the measure is symmetrized by taking $\max(M(A, B), M(B, A))$.

Comparison of Interestingness Measures

- Null-(transaction) invariance is crucial for correlation analysis
- Lift and χ^2 are not null-invariant
- 5 null-invariant measures

	Milk	No Milk	Sum (row)
Coffee	m, c	$\sim m, c$	c
No Coffee	m, $\sim c$	$\sim m, \sim c$	$\sim c$
Sum (col.)	m	$\sim m$	Σ

Measure	Definition	Range	Null-Invariant
$\chi^2(a, b)$	$\sum_{i,j=0,1} \frac{(e(a_i, b_j) - o(a_i, b_j))^2}{e(a_i, b_j)}$	[0, ∞]	No
Lift(a, b)	$\frac{P(ab)}{P(a)P(b)}$	[0, ∞]	No
AllConf(a, b)	$\frac{sup(ab)}{\max\{sup(a), sup(b)\}}$	[0, 1]	Yes
Coherence(a, b)	$\frac{sup(ab)}{sup(a)+sup(b)-sup(ab)}$	[0, 1]	Yes
Cosine(a, b)	$\frac{sup(ab)}{\sqrt{sup(a)sup(b)}}$	[0, 1]	Yes
Kulc(a, b)	$\frac{sup(ab)}{2} (\frac{1}{sup(a)} + \frac{1}{sup(b)})$	[0, 1]	Yes
MaxConf(a, b)	$\max\{\frac{sup(ab)}{sup(a)}, \frac{sup(ab)}{sup(b)}\}$	[0, 1]	Yes

Table 3. Interestingness measure definitions.

Null-transactions w.r.t. m and c

Kulczynski measure (1927)

Null-invariant

Data set	mc	$\bar{m}\bar{c}$	$m\bar{c}$	$\bar{m}c$	χ^2	Lift	AllConf	Coherence	Cosine	Kulc	MaxConf
D_1	10,000	1,000	1,000	100,000	90557	9.26	0.91	0.83	0.91	0.91	0.91
D_2	10,000	1,000	1,000	100	0	1	0.91	0.83	0.91	0.91	0.91
D_3	100	1,000	1,000	100,000	670	8.44	0.09	0.05	0.09	0.09	0.09
D_4	1,000	1,000	1,000	100,000	24740	25.75	0.5	0.33	0.5	0.5	0.5
D_5	1,000	100	10,000	100,000	8173	9.18	0.09	0.09	0.29	0.5	0.91
D_6	1,000	10	100,000	100,000	965	1.97	0.01	0.01	0.10	0.5	0.99

Table 2. Example data sets. Subtle: They disagree

Summary

- ◆ Basic concepts: association rules, support-confident framework, closed and max-patterns
- ◆ Scalable frequent pattern mining methods
 - Apriori (Candidate generation & test)
 - Projection-based (FPgrowth, CLOSET+, ...)
 - Vertical format approach (ECLAT, CHARM, ...)
- ◆ Which patterns are interesting?
 - Pattern evaluation methods

42

Ref: Basic Concepts of Frequent Pattern Mining

- ◆ (Association Rules) R. Agrawal, T. Imielinski, and A. Swami. Mining association rules between sets of items in large databases. SIGMOD'93.
- ◆ (Max-pattern) R. J. Bayardo. Efficiently mining long patterns from databases. SIGMOD'98.
- ◆ (Closed-pattern) N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Discovering frequent closed itemsets for association rules. ICDT'99.
- ◆ (Sequential pattern) R. Agrawal and R. Srikant. Mining sequential patterns. ICDE'95

43

Ref: Apriori and Its Improvements

- ◆ R. Agrawal and R. Srikant. Fast algorithms for mining association rules. VLDB'94.
- ◆ H. Mannila, H. Toivonen, and A. I. Verkamo. Efficient algorithms for discovering association rules. KDD'94.
- ◆ A. Savasere, E. Omiecinski, and S. Navathe. An efficient algorithm for mining association rules in large databases. VLDB'95.
- ◆ J. S. Park, M. S. Chen, and P. S. Yu. An effective hash-based algorithm for mining association rules. SIGMOD'95.
- ◆ H. Toivonen. Sampling large databases for association rules. VLDB'96.
- ◆ S. Brin, R. Motwani, J. D. Ullman, and S. Tsur. Dynamic itemset counting and implication rules for market basket analysis. SIGMOD'97.
- ◆ S. Sarawagi, S. Thomas, and R. Agrawal. Integrating association rule mining with relational database systems: Alternatives and implications. SIGMOD'98.

Ref: Depth-First, Projection-Based FP Mining

- ◆ R. Agarwal, C. Aggarwal, and V. V. V. Prasad. A tree projection algorithm for generation of frequent itemsets. J. Parallel and Distributed Computing:02.
- ◆ J. Han, J. Pei, and Y. Yin. Mining frequent patterns without candidate generation. SIGMOD'00.
- ◆ J. Liu, Y. Pan, K. Wang, and J. Han. Mining Frequent Item Sets by Opportunistic Projection. KDD'02.
- ◆ J. Han, J. Wang, Y. Lu, and P. Tzvetkov. Mining Top-K Frequent Closed Patterns without Minimum Support. ICDM'02.
- ◆ J. Wang, J. Han, and J. Pei. CLOSET+: Searching for the Best Strategies for Mining Frequent Closed Itemsets. KDD'03.
- ◆ G. Liu, H. Lu, W. Lou, J. X. Yu. On Computing, Storing and Querying Frequent Patterns. KDD'03.
- ◆ G. Grahne and J. Zhu, Efficiently Using Prefix-Trees in Mining Frequent Itemsets, Proc. ICDM'03 Int. Workshop on Frequent Itemset Mining Implementations (FIMI'03), Melbourne, FL, Nov. 2003

Ref: Vertical Format and Row Enumeration Methods

- ◆ M. J. Zaki, S. Parthasarathy, M. Ogihara, and W. Li. Parallel algorithm for discovery of association rules. DAMI:97.
- ◆ Zaki and Hsiao. CHARM: An Efficient Algorithm for Closed Itemset Mining, SDM'02.
- ◆ C. Bucila, J. Gehrke, D. Kifer, and W. White. DualMiner: A Dual-Pruning Algorithm for Itemsets with Constraints. KDD'02.
- ◆ F. Pan, G. Cong, A. K. H. Tung, J. Yang, and M. Zaki, CARPENTER: Finding Closed Patterns in Long Biological Datasets. KDD'03.
- ◆ H. Liu, J. Han, D. Xin, and Z. Shao, Mining Interesting Patterns from Very High Dimensional Data: A Top-Down Row Enumeration Approach, SDM'06.

46

Ref: Mining Correlations and Interesting Rules

- ◆ M. Klemettinen, H. Mannila, P. Ronkainen, H. Toivonen, and A. I. Verkamo. Finding interesting rules from large sets of discovered association rules. CIKM'94.
- ◆ S. Brin, R. Motwani, and C. Silverstein. Beyond market basket: Generalizing association rules to correlations. SIGMOD'97.
- ◆ C. Silverstein, S. Brin, R. Motwani, and J. Ullman. Scalable techniques for mining causal structures. VLDB'98.
- ◆ P.-N. Tan, V. Kumar, and J. Srivastava. Selecting the Right Interestingness Measure for Association Patterns. KDD'02.
- ◆ E. Omiecinski. Alternative Interest Measures for Mining Associations. TKDE'03.
- ◆ T. Wu, Y. Chen and J. Han, "Association Mining in Large Databases: A Re-Examination of Its Measures", PKDD'07

47

Ref: Freq. Pattern Mining Applications

- ◆ Y. Huhtala, J. Kärkkäinen, P. Porkka, H. Toivonen. Efficient Discovery of Functional and Approximate Dependencies Using Partitions. ICDE'98.
- ◆ H. V. Jagadish, J. Madar, and R. Ng. Semantic Compression and Pattern Extraction with Fascicles. VLDB'99.
- ◆ T. Dasu, T. Johnson, S. Muthukrishnan, and V. Shkapenyuk. Mining Database Structure; or How to Build a Data Quality Browser. SIGMOD'02.
- ◆ K. Wang, S. Zhou, J. Han. Profit Mining: From Patterns to Actions. EDBT'02.

48



Questions?