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

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: One man's trash could be another man's treasure
 - □ Query-based: Relevant to a user's particular request
 - □ Against one's knowledge-base: unexpected, freshness, timeliness
 - Visualization tools: Multi-dimensional, interactive examination

Limitation of the Support-Confidence Framework

- □ Are *s* and *c* interesting in association rules: "A \Rightarrow B" [*s*, *c*]? Be careful!
- 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 2-	Way Cont
not eat-cereal	200	50	250	^{way contingency table}
sum(col.)	600	400	1000	

- □ Association rule mining may generate the following:
 - \Box 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:
 - □ ¬ play-basketball \Rightarrow eat-cereal [35%, 87.5%] (high s & c)

Interestingness Measure: Lift

□ Measure of dependent/correlated events: lift $lift(B,C) = \frac{c(B \rightarrow C)}{s(C)} = \frac{s(B \cup C)}{s(B) \times s(C)}$

□ Lift(B, C) may tell how B and C are correlated

□ Lift(B, C) = 1: B and C are independent

- \Box > 1: positively correlated
- □ < 1: negatively correlated

□ For 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;

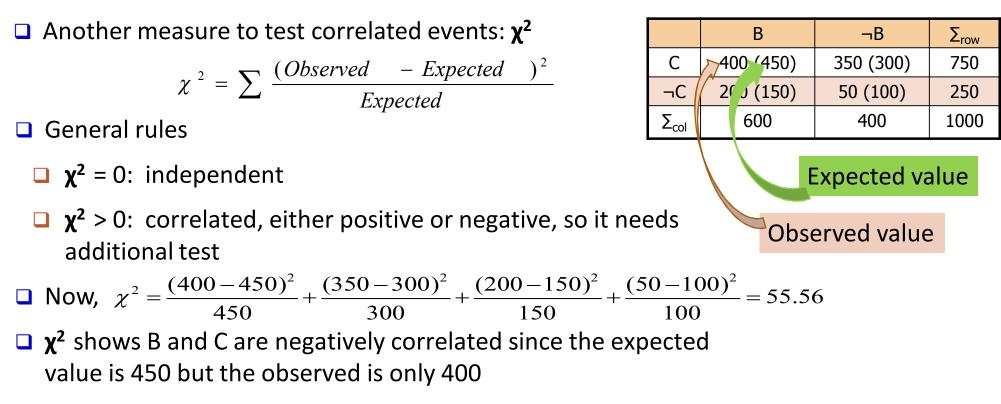
B and ¬C are positively correlated since lift(B, ¬C) > 1

Lift is more telling than s & c

	В	¬Β	Σ _{row}
С	400	350	750
٦C	200	00 50 250	
$\Sigma_{col.}$	600	400	1000

33

Interestingness Measure: χ^2



 $\hfill \chi^2$ is also more telling than the support-confidence framework

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

- Null transactions: Transactions that contain neither B nor C
- Let's examine the 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!)
- $\Box \chi^2 = 670$: Observed(BC) >> expected value (11.85)
- Too many null transactions may "spoil the soup"!

	В	⊐B	Σ _{row}				
C	100	1000	1100				
٦C	1000	100000	101000				
$\Sigma_{col.}$	1100 🧯	101000	102100				
	null transactions						

Contingency table with expected values added

	В	⊐B	Σ _{row}
С	100 (11.85)	1000	1100
٦C	1000 (988.15)	100000	101000
$\Sigma_{col.}$	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]
$\chi^2(A,B)$	$\sum_{i,j=0,1} \frac{(e(a_i b_j) - o(a_i b_j))^2}{e(a_i b_j)}$	$[0,\infty]$	No	X² and lift are not
Lift(A, B)	$\frac{s(A \cup B)}{s(A) \times s(B)}$	$[0,\infty]$	No	null-invariant
AllConf(A, B)	$\frac{s(A \cup B)}{max\{s(A), s(B)\}}$	[0,1]	Yes	laccard concine
Jaccard(A, B)	$\frac{s(A \cup B)}{s(A) + s(B) - s(A \cup B)}$	[0,1]	Yes	Jaccard, consine, AllConf, MaxConf,
Cosine(A, B)	$\frac{s(A \cup B)}{\sqrt{s(A) \times s(B)}}$	[0,1]	Yes	and Kulczynski
Kulczynski(A, B)	$\frac{\frac{1}{2}\left(\frac{s(A\cup B)}{s(A)} + \frac{s(A\cup B)}{s(B)}\right)}{\frac{1}{2}\left(\frac{s(A\cup B)}{s(B)} + \frac{s(A\cup B)}{s(B)}\right)}$	[0,1]	Yes	are null-invariant measures
MaxConf(A, B)	$max\{\frac{s(A)}{s(A\cup B)}, \frac{s(B)}{s(A\cup B)}\}$	[0,1]	Yes	

Null Invariance: An Important Property

- U 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}	m	$\neg m$	Σ

- milk vs. coffee contingency table
- Lift and χ² 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 w.r.t. m and c

Data set	mc	$\neg mc$	$m\neg c$	$m\neg c$	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?
 - \Box D₄-D₆ 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}
coffee	mc	$\neg mc$	c
$\neg coffee$	$m \neg c$	$\neg m \neg c$	$\neg c$
Σ_{col}	m	$\neg m$	Σ

					<u> </u>				
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	100,000	0.01	0.01	0.10	0.5	0.99

All 5 are null-invariant

Subtle: They disagree on those cases

Analysis of DBLP Coauthor Relationships

Recent DB conferences,	removing balanc	ed associations,	low sup, etc.
	Territe Ting Balance		

ID	Author A	Author B	$s(A \cup B)$	s(A)	s(B)	Jaccard	Cosine	Kulc
1	Hans-Peter Kriegel	Martin Ester	28	146	54	0.163(2)	0.315(7)	0.355(9)
2	Michael Carey	Miron Livny	26	104	<mark>58</mark>	0.191(1)	0.335(4)	0.349(10)
3	Hans-Peter Kriegel	Joerg Sander	24	146	36	0.152(3)	0.331(5)	0.416(8)
4	Christos Faloutsos	Spiros Papadimitriou	20	162	26	0.119(7)	0.308(10)	0.446(7)
5	Hans-Peter Kriegel	Martin Pfeifle	18	146	18	0.123~(6)	0.351(2)	0.562(2)
6	Hector Garcia-Molina	Wilburt Labio	16	144	18	0.110(9)	0.314(8)	0.500(4)
7	Divyakant Agrawal	Wang Hsiung	16	120	16	0.133(5)	0.365(1)	0.567(1)
8	Elke Rundensteiner	Murali Mani	16	104	20	0.148(4)	0.351(3)	0.477(6)
9	Divyakant Agrawal	Oliver Po	\bigtriangledown 12	120	12	0.100(10)	0.316(6)	0.550(3)
10	Gerhard Weikum	Martin Theobald	12	106	14	0.111(8)	0.312(9)	0.485(5)
22. 93.				20	s Dh			

Advisor-advisee relation: Kulc: high, Jaccard: low, cosine: middle

- Which pairs of authors are strongly related?
 - Use Kulc to find Advisor-advisee, close collaborators

39

Imbalance Ratio with Kulczynski Measure

IR (Imbalance Ratio): measure the imbalance of two itemsets A and B in rule implications: 1. (1) (D)

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

- Kulczynski and Imbalance Ratio (IR) together present a clear picture for all the three datasets D_4 through D_6
 - D_4 is neutral & balanced; D_5 is neutral but imbalanced
 - 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	$\bigcirc 0.5$	0.5	0
D_5	1,000	100	10,000	100,000	0.09	0.29	0.5	0.89
D_6	1,000	10	100,000	100,000	0.01	0.10	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
- \Box Lift, χ^2 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
- Exercise: Mining research collaborations from research bibliographic data
 - □ Find a group of frequent collaborators from research bibliographic data (e.g., DBLP)
 - Can you find the likely advisor-advisee relationship and during which years such a relationship happened?
 - Ref.: C. Wang, J. Han, Y. Jia, J. Tang, D. Zhang, Y. Yu, and J. Guo, "Mining Advisor-Advisee Relationships from Research Publication Networks", KDD'10

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



Summary: Mining Frequent Patterns, Association and Correlations

- Basic Concepts:
 - Frequent Patterns, Association Rules, Closed Patterns and Max-Patterns
- Frequent Itemset Mining Methods
 - The Downward Closure Property and 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
- Which Patterns Are Interesting?—Pattern Evaluation Methods
- \Box Interestingness Measures: Lift and χ^2
- Null-Invariant Measures
- Comparison of Interestingness Measures