

Data Mining:

Concepts and Techniques

(3rd ed.)

slightly adapted slides

— Chapter 6 —

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

©2013 - 2018 Han, Kamber & Pei. All rights reserved.

29

Scalable Frequent Itemset 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

30

Basic Concepts: Frequent Patterns

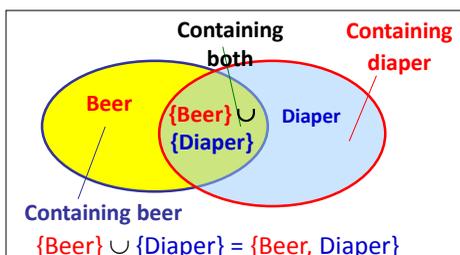
Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

- **itemset**: A set of one or more items
 - **k-itemset** $X = \{x_1, \dots, x_k\}$
 - **(absolute) support (count)** of X : Frequency or the number of occurrences of an itemset X
 - **(relative) support, s** : The fraction of transactions that contains X (i.e., the probability that a transaction contains X)
 - An itemset X is **frequent** if the support of X is no less than a **minsup** threshold (σ)
- Let **minsup** = 50%
 - Freq. 1-itemsets:
 - Beer: 3 (60%); Nuts: 3 (60%)
 - Diaper: 4 (80%); Eggs: 3 (60%)
 - Freq. 2-itemsets:
 - {Beer, Diaper}: 3 (60%)

31

From Frequent Itemsets to 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



Note: Itemset: $X \cup Y$, a subtle notation!

- **Association rules**: $X \rightarrow Y (s, c)$
 - **Support, s** : The probability that a transaction contains $X \cup Y$
 - **Confidence, c** : The conditional probability that a transaction containing X also contains Y
 - $c = \text{sup}(X \cup Y) / \text{sup}(X)$
- **Association rule mining**: Find **all** of the rules, $X \rightarrow Y$, with minimum support and confidence
- **Frequent itemsets**: Let **minsup** = 50%
 - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
 - Freq. 2-itemsets: {Beer, Diaper}: 3
- **Association rules**: Let **minconf** = 50%
 - $Beer \rightarrow Diaper$ (60%, 100%)
 - $Diaper \rightarrow Beer$ (60%, 75%)

32

Frequent Itemset Mining

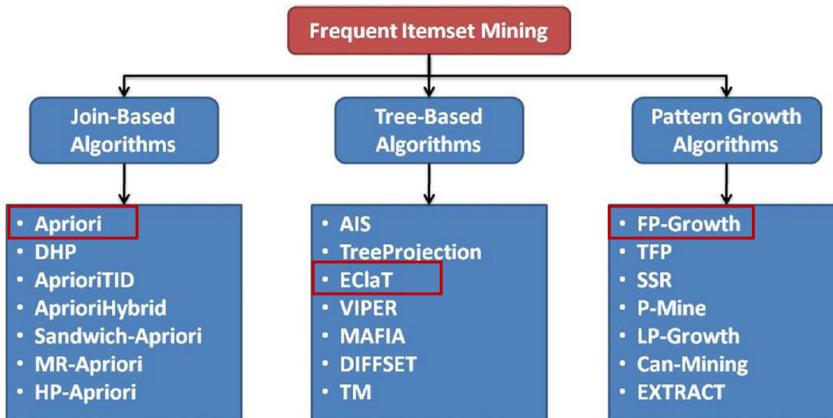


Fig. 12 Classification of Frequent Pattern Mining algorithms

Figure from: C. Chin-Hoong et al., Algorithms for frequent itemset mining: a literature review", Artificial Intelligence Review, Springer, 2018

November 8, 2018

Data Mining: Concepts and Techniques

33

FPGrowth: Mining Frequent Patterns by Pattern Growth

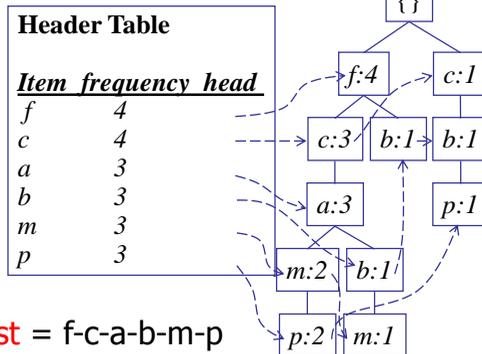
- Idea: Frequent pattern growth (FPGrowth)
 - Find frequent single items and partition the database based on each such item
 - Recursively grow frequent patterns by doing the above for each partitioned database (also called *conditional database*)
 - To facilitate efficient processing, an efficient data 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

Construct FP-tree from a Transaction Database

TID	Items bought	(ordered) frequent items
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}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

min_support = 3

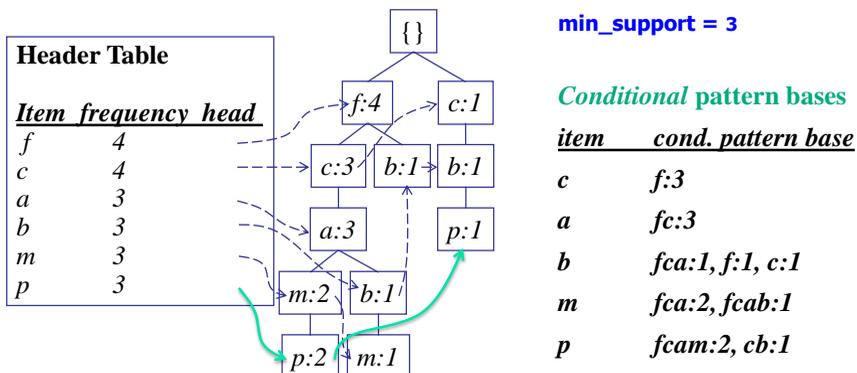
1. Scan DB once, find frequent 1-itemset (single item pattern)
2. Sort frequent items in frequency descending order, **f-list**
3. Scan DB again, construct FP-tree



35

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

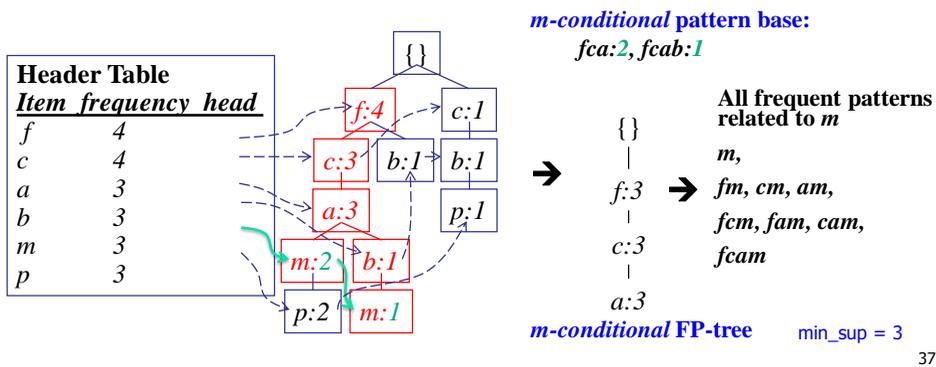
- Pattern mining can be partitioned according to current patterns
 - Patterns containing **p**: **p**'s conditional database: *fcam:2, cb:1*
 - Patterns having **m** but no **p**: **m**'s conditional database: *fca:2, fcab:1*
 -
- p**'s conditional pattern base: *transformed prefix paths* of item **p**



36

FP-tree Mining: From Conditional Pattern-bases to Conditional FP-trees

- For each conditional pattern-base
 - Accumulate the count for each item in the base
 - Construct the conditional FP-tree for the frequent items of the conditional pattern base



37

Mine Each Conditional Pattern-Base Recursively

Conditional pattern bases

item base	cond. pattern
	min_support = 3
c	f:3
a	fc:3
b	fca:1, f:1, c:1
m	fca:2, fcab:1
p	fcam:2, cb:1

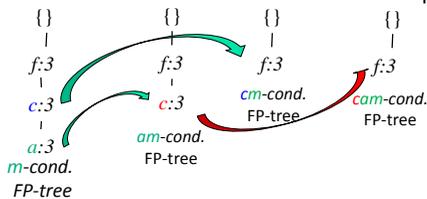
- For each conditional pattern-base
 - Mine single-item patterns *c*, *a*, *b*, ...
 - Construct its FP-tree & mine it recursively

p-conditional PB: fcam:2, cb:1 → c: 3

m-conditional PB: fca:2, fcab:1 → fca: 3

b-conditional PB: fca:1, f:1, c:1 → φ

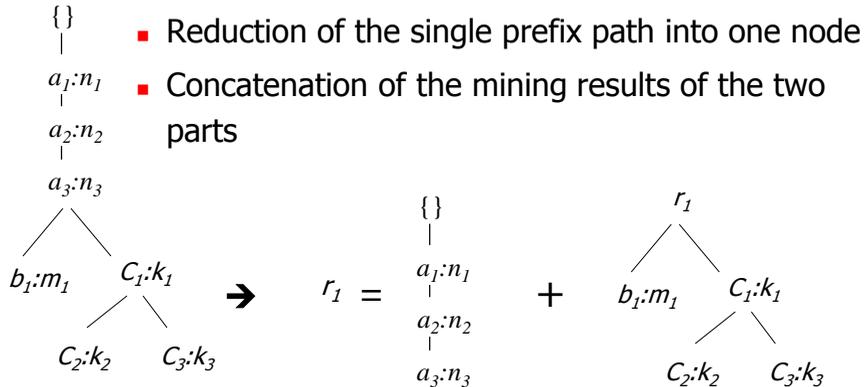
Actually, for single branch FP-tree, all frequent patterns can be generated in one shot



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



39

The Apriori Algorithm—An Example

minsup = 2

Database TDB

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

1st scan

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

L_1

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

L_2

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

C_2

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

2nd scan

Itemset
{A, B}
{A, C}
{A, E}
{B, C}
{B, E}
{C, E}

C_3

Itemset
{B, C, E}

3rd scan

Itemset	sup
{B, C, E}	2

L_3

40

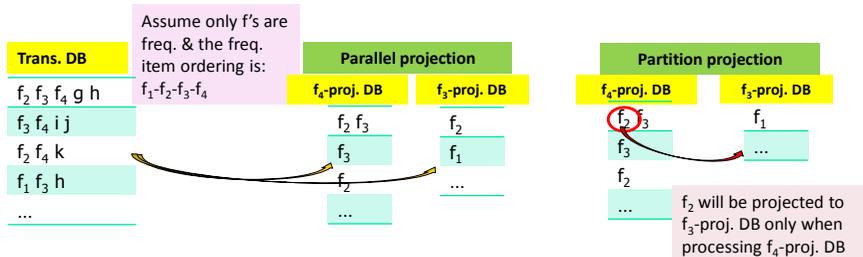
Benefits of the FP-tree Structure

- **Completeness**
 - Preserve complete information for frequent pattern mining
 - Never break a long pattern of any transaction
- **Compactness**
 - *Reduce irrelevant info* —infrequent items are gone
 - *Items in frequency descending order*: the more frequently occurring, the more likely to be shared
 - *Never be larger than the original database* (if not counting: node-links and the *count* field)

41

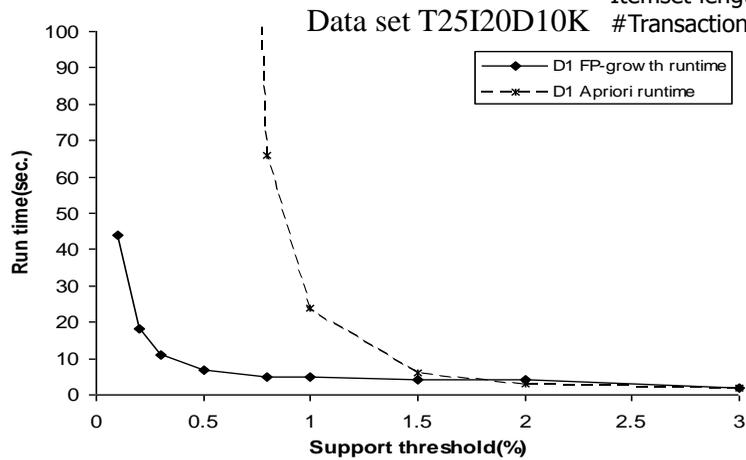
Scaling FP-growth by Database Projection

- What if FP-tree cannot fit in memory? — DB projection
 - Project the DB based on patterns
 - 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



FP-Growth vs. Apriori: Scalability With the Support Threshold

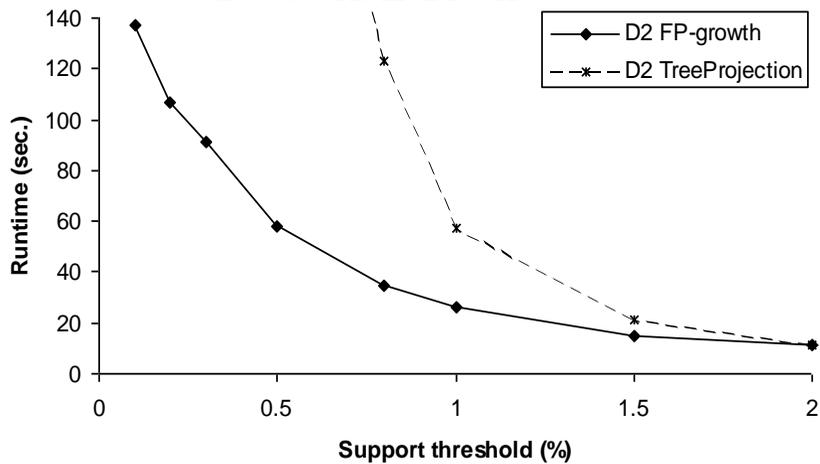
IBM Synth.: Average
Transaction length: 25
Itemset length: 10K
#Transactions: 20K



43

FP-Growth vs. Tree-Projection: Scalability with the Support Threshold

Data set T25I20D100K



44

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 operations: 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)

45

Extension of Pattern Growth Mining Methodology

- Mining closed frequent itemsets and max-patterns
 - CLOSET (DMKD'00), FPclose, and FPMMax (Grahne & Zhu, Fimi'03)
- Mining sequential patterns
 - PrefixSpan (ICDE'01), CloSpan (SDM'03), BIDE (ICDE'04)
- Mining graph patterns
 - gSpan (ICDM'02), CloseGraph (KDD'03)
- Constraint-based mining of frequent patterns
 - Convertible constraints (ICDE'01), gPrune (PAKDD'03)
- Computing iceberg data cubes with complex measures
 - H-tree, H-cubing, and Star-cubing (SIGMOD'01, VLDB'03)
- Pattern-growth-based Clustering
 - MaPle (Pei, et al., ICDM'03)
- Pattern-Growth-Based Classification
 - Mining frequent and discriminative patterns (Cheng, et al, ICDE'07)

46

Scalable Frequent Itemset 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



47

Closed Patterns and Max-Patterns

An itemset X is a **closed pattern**

- if X is *frequent* and
- there exists *no super-pattern* $Y \supset X$, with the same support as X

An itemset X is a **max-pattern**

- if X is frequent and
- there exists no frequent super-pattern $Y \supset X$

A Closed pattern is a lossless compression of freq. patterns

- Reducing the # of patterns and rules

48

CLOSET+: Mining Closed Itemsets by Pattern-Growth

- Efficient, *direct* mining of *closed* itemsets
- Ex. Itemset merging: If *Y* appears in every occurrence of *X*, then *Y* is merged with *X*
 - d-proj. db: {*ac*ef, *ac*f} → *acfd*-proj. db: {*e*}
 - thus we get: *acfd*:2
- Many other tricks (but not detailed here), such as
 - Hybrid tree projection
 - Bottom-up physical tree-projection
 - Top-down pseudo tree-projection
 - Sub-itemset pruning
 - Item skipping
 - Efficient subset checking
- For details, see J. Wang, et al., "CLOSET+:", KDD'03

TID	Items
1	acdef
2	abe
3	cefg
4	acdf

Let minsupport = 2

a:3, c:3, d:2, e:3, f:3

F-List: a-c-e-f-d

MaxMiner: Mining Max-Patterns

- 1st scan: find frequent items
 - A, B, C, D, E
- 2nd scan: find support for
 - AB, AC, AD, AE, **ABCDE**
 - BC, BD, BE, **BCDE**
 - CD, CE, **CDE**
 - DE
- Since **BCDE** is a max-pattern, no need to check **BCD**, **BDE**, **CDE** in later scan
- R. Bayardo. *Efficiently mining long patterns from databases.* *SIGMOD'98*

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

minsup = 2

Potential
max-patterns

A MapReduce-Based Parallel Frequent Pattern Growth Algorithm for Spatiotemporal Association Analysis of Mobile Trajectory Big Data

Dawen Xia ^{1,2}, Xiaonan Lu,¹ Huaqing Li ³, Wendong Wang,² Yantao Li ², and Zili Zhang^{2,4}

¹College of Data Science and Information Engineering and College of National Culture and Cognitive Science, Guizhou Minzu University, Guiyang 550025, China

²College of Computer and Information Science, Southwest University, Chongqing 400715, China

³College of Electronic and Information Engineering, Southwest University, Chongqing 400715, China

⁴School of Information Technology, Deakin University, Geelong, VIC 3216, Australia

Correspondence should be addressed to Dawen Xia; dwxia@gzmu.edu.cn and Huaqing Li; huaqingli@hotmail.com

Received 14 June 2017; Accepted 13 November 2017; Published 28 January 2018

Academic Editor: Michele Scarpiniti

Copyright © 2018 Dawen Xia et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Frequent pattern mining is an effective approach for spatiotemporal association analysis of mobile trajectory big data in data-driven intelligent transportation systems. While existing parallel algorithms have been successfully applied to frequent pattern mining of large-scale trajectory data, two major challenges are how to overcome the inherent defects of Hadoop to cope with taxi trajectory big data including massive small files and how to discover the implicitly spatiotemporal frequent patterns with MapReduce. To conquer these challenges, this paper presents a MapReduce-based Parallel Frequent Pattern growth (MR-PPF) algorithm to analyze the spatiotemporal characteristics of taxi operating using large-scale taxi trajectories with massive small file processing strategies on a Hadoop platform. More specifically, we first implement three methods, that is, Hadoop Archives (HAR), CombineFileInputFormat (CFIF), and Sequence Files (SF), to overcome the existing defects of Hadoop and then propose two strategies based on their performance evaluations. Next, we incorporate SF into Frequent Pattern growth (FP-growth) algorithm and then implement the optimized FP-growth algorithm on a MapReduce framework. Finally, we analyze the characteristics of taxi operating in both spatial and temporal dimensions by MR-PPF in parallel. The results demonstrate that MR-PPF is superior to existing Parallel FP-growth (PPF) algorithm in efficiency and scalability.

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**

Interestingness: Limitation of the Support-Confidence Framework

Be careful!

- 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 (total)

2-way contingency table

- 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**

s = support
c = confidence

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

Lift is more telling than *s* & *c*

	B	¬B	Σ _{row}
C	400	350	750
¬C	200	50	250
Σ _{col}	600	400	1000 (total)

B = Play Basketball
C = Eat Cereal

- Lift(B, C) may tell how B and C are correlated
 - Lift(B, C) = 1: B and C are independent
 - > 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

Interestingness Measure: χ^2

- Another measure to test correlated events: χ^2

$$\chi^2 = \sum \frac{(Observedcount - Expectedcount)^2}{Expectedcount}$$

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

- General rules

- $\chi^2 = 0$: independent
- $\chi^2 > 0$: correlated, either positive or negative, so it needs additional test to determine which correlation

- Now,

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

- χ^2 shows B and C are negatively correlated since the expected value is 450 (= 600 * 750/1000) but the observed is lower, only 400
- χ^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)
- In these transactions it is unlikely that B & C will happen together!

- But, $Lift(B, C) = 8.44 \gg 1$

(Lift shows B and C are strongly positively correlated!)

- $\chi^2 = 670$: Observed(BC) >> expected value (11.85)

- Too many null transactions may "spoil the soup"!

dataset D

	B	-B	Σ_{row}
C	100	1000	1100
-C	1000	100000	101000
$\Sigma_{col.}$	1100	101000	102100

null transactions

Contingency table with expected values added

	B	-B	Σ_{row}
C	100 (11.85)	1000	1100
-C	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
$Lift(A, B)$	$\frac{s(A \cup B)}{s(A) \times s(B)}$	$[0, \infty]$	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
$Cosine(A, B)$	$\frac{s(A \cup B)}{\sqrt{s(A) \times s(B)}}$	$[0, 1]$	Yes
$Kulczynski(A, B)$	$\frac{1}{2} \left(\frac{s(A \cup B)}{s(A)} + \frac{s(A \cup B)}{s(B)} \right)$	$[0, 1]$	Yes
$MaxConf(A, B)$	$\max\left\{ \frac{s(A)}{s(A \cup B)}, \frac{s(B)}{s(A \cup B)} \right\}$	$[0, 1]$	Yes

χ^2 and lift are not null-invariant

Jaccard, Cosine, AllConf, MaxConf, and Kulczynski are null-invariant measures

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 vs. coffee 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

- 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 w.r.t. m and c

Data set	mc	$\neg mc$	$m\neg c$	$\neg m\neg c$	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

Assignment: Check the interestingness measures in the table.

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

All 5 are null-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	100,000	0.91	0.01	0.10	0.5	0.99

Subtle: They disagree on those cases

Analysis of DBLP Coauthor Relationships

Recent DB conferences, removing balanced associations, low sup, etc.

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	58	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	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)

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

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

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 D_4 through D_6
 - D_4 is neutral & balanced
 - D_5 is neutral but imbalanced
 - D_6 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	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?

Optional Reading: 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
- 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, 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

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
 - Interestingness Measures: Lift and χ^2
 - Null-Invariant Measures
 - Comparison of Interestingness Measures

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

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.

67

References (II) Efficient Pattern Mining Methods

- R. Agrawal and R. Srikant, "Fast algorithms for mining association rules", VLDB'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
- S. Sarawagi, S. Thomas, and R. Agrawal, "Integrating association rule mining with relational database systems: Alternatives and implications", SIGMOD'98
- M. J. Zaki, S. Parthasarathy, M. Ogihara, and W. Li, "Parallel algorithm for discovery of association rules", Data Mining and Knowledge Discovery, 1997
- J. Han, J. Pei, and Y. Yin, "Mining frequent patterns without candidate generation", SIGMOD'00
- M. J. Zaki and Hsiao, "CHARM: An Efficient Algorithm for Closed Itemset Mining", SDM'02
- J. Wang, J. Han, and J. Pei, "CLOSET+: Searching for the Best Strategies for Mining Frequent Closed Itemsets", KDD'03
- C. C. Aggarwal, M.A., Bhuiyan, M. A. Hasan, "Frequent Pattern Mining Algorithms: A Survey", in Aggarwal and Han (eds.): Frequent Pattern Mining, Springer, 2014

References (III) Pattern Evaluation

- C. C. Aggarwal and P. S. Yu. A New Framework for Itemset Generation. PODS'98
- S. Brin, R. Motwani, and C. Silverstein. Beyond market basket: Generalizing association rules to correlations. SIGMOD'97
- M. Klemettinen, H. Mannila, P. Ronkainen, H. Toivonen, and A. I. Verkamo. Finding interesting rules from large sets of discovered association rules. CIKM'94
- E. Omiecinski. Alternative Interest Measures for Mining Associations. TKDE'03
- P.-N. Tan, V. Kumar, and J. Srivastava. Selecting the Right Interestingness Measure for Association Patterns. KDD'02
- T. Wu, Y. Chen and J. Han, Re-Examination of Interestingness Measures in Pattern Mining: A Unified Framework, Data Mining and Knowledge Discovery, 21(3):371-397, 2010