



CS 412 Intro. to Data Mining

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

Jiawei Han, Computer Science, Univ. Illinois at Urbana-Champaign, 2017

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

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



mine patterns in various Data



What Is Pattern Discovery?

การค้นหากราฟข้อมูล

- What are patterns? *กําหนดค่าของสิ่งของที่มีรูปแบบ (patterns ตรงๆ ของ dataset), set of items จํากัดในช่วงเวลา*
- 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 *ตัวอย่างที่ค้นพบในชุดข้อมูล*
- Motivation examples:
 - What products were often purchased together? *สินค้าที่ซื้อกันบ่อยๆ*
 - What are the subsequent purchases after buying an iPad? *หลังจากซื้อ iPad ไปแล้ว ผู้คนมักซื้อสิ่งของอื่นๆ อีก*
 - What code segments likely contain copy-and-paste bugs? *ชิ้นโค้ด ซึ่งมีลักษณะคล้ายกันมาก หรือมี syntax ที่ซ้ำๆ กัน*
 - What word sequences likely form phrases in this corpus? *คำศัพท์ที่มักจะเป็นวลีในชุดข้อมูลนี้*

Basic Concepts: k-Itemsets and Their Supports

សំណើអាជ្ញាធរកំណត់

មេដាច់ទិន្នន័យទាំងអស់ជាអំពី ទំនួរ

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

មេដាច់ទិន្នន័យ និង ការបង្កើត

- **k-itemset:** $X = \{x_1, \dots, x_k\}$

ការបង្កើតទិន្នន័យ និង ការបង្កើត 3 ចំណាំ

- Ex. {Beer, Nuts, Diaper} is a 3-itemset

តុលាភាសា transaction និង support និង ការបង្កើត និង ការគិតថ្មី

- **(absolute) support (count)** of X, $\text{sup}\{X\}$:

ការបង្កើត

Frequency or the number of occurrences of an itemset X

- Ex. $\text{sup}\{\text{Beer}\} = 3$ តុលាភាសា transaction និង ការបង្កើត

- Ex. $\text{sup}\{\text{Diaper}\} = 4$ តុលាភាសា transaction និង ការបង្កើត

- Ex. $\text{sup}\{\text{Beer}, \text{Diaper}\} = 3$ តុលាភាសា transaction និង ការបង្កើត

- Ex. $\text{sup}\{\text{Beer}, \text{Eggs}\} = 1$

Tid	transaction ID	Items bought
10	transaction 1	Beer, Nuts, Diaper
20	transaction 2	Beer, Coffee, Diaper
30	transaction 3	Beer, Diaper, Eggs
40	transaction 4	Nuts, Eggs, Milk
50	transaction 5	Nuts, Coffee, Diaper, Eggs, Milk

តុលាភាសា transaction និង support និង ការបង្កើត

- **(relative) support**, $s\{X\}$: The fraction of transactions that contains X (i.e., the probability that a transaction contains X)

- Ex. $s\{\text{Beer}\} = 3/5 = 60\%$ តុលាភាសា transaction និង Beer

- Ex. $s\{\text{Diaper}\} = 4/5 = 80\%$

- Ex. $s\{\text{Beer}, \text{Eggs}\} = 1/5 = 20\%$

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
- Let $\sigma = 50\%$ (σ : *minsup threshold*)

For the given 5-transaction dataset

- All the frequent 1-itemsets:
 - Beer: 3/5 (60%); Nuts: 3/5 (60%)
 - Diaper: 4/5 (80%); Eggs: 3/5 (60%)
- All the frequent 2-itemsets:
 - {Beer, Diaper}: 3/5 (60%)
- All the frequent 3-itemsets?
- None

• Coffee : 2/5 (40%) = ໃຫຍ່ນໜີ້ນີ້

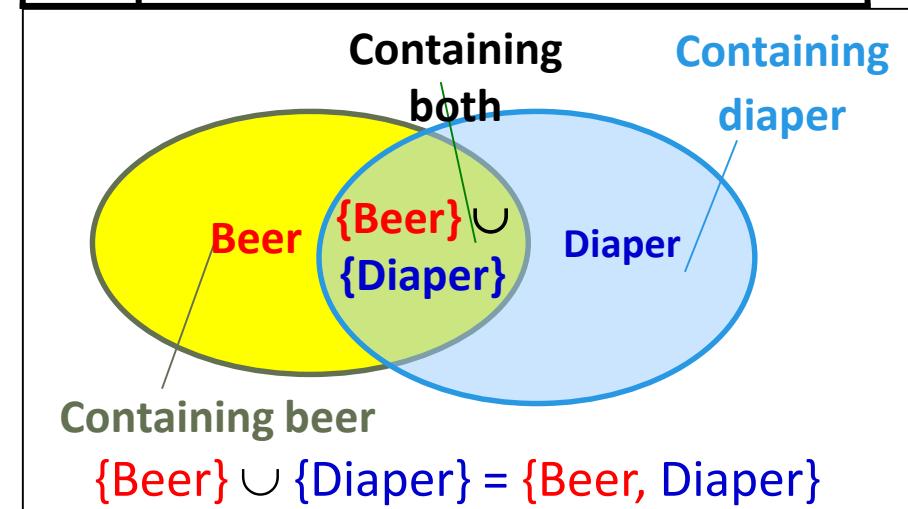
Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

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

From Frequent Itemsets to Association Rules

- Comparing with itemsets, rules can be more telling
 - Ex. $\text{Diaper} \rightarrow \text{Beer}$
- *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. $s\{\text{Diaper}, \text{Beer}\} = 3/5 = 0.6$ (i.e., 60%)
 - Confidence, c: The conditional probability that a transaction containing X also contains Y
 - Calculation: $c = \frac{s(X \cup Y)}{s(X)} = \frac{0.6}{0.5} = 1.2$
 - Ex. $c = \frac{s\{\text{Diaper}, \text{Beer}\}}{s\{\text{Diaper}\}} = \frac{3/5}{5/5} = 0.75$

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



Note: $X \cup Y$: the union of two itemsets

- The set contains both X and Y

Mining Frequent Itemsets and Association Rules

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

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

(Q: Are these all rules?)

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
- FP-Growth: A Frequent Pattern-Growth Approach
- Mining Closed Patterns

Apriori Pruning and Scalable Mining Methods

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

Apriori: A Candidate Generation & Test Approach

ଅପ୍ରିଓରି କାନ୍ଡିଡେଟ ଗେନେରେସନ୍ ଏବଂ ଟେସ୍ଟ ଅପ୍ରାକ୍ରମ

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

The Apriori Algorithm—An Example

Database TDB

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

↳ 1 transaction

กำหนด minsup
minsup = 2

C_1

1st scan
ใน 1 รายการ

มอง One Itemset

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

F_1

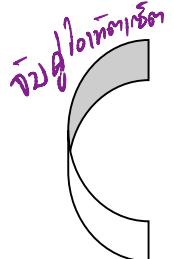
จัดเรียงตาม Itemset ค่า sup มากไปน้อย

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

จัดเรียงตาม sup

F_2

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



C_2

จัดเรียง

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

C_2

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



2nd scan

ค่า sup ของ itemset ที่มี

3 รายการ

Itemset
{B, C, E}

3rd scan

Itemset	sup
{B, C, E}	2

สร้างเป็น tree Itemset
ตามจริง ใจ 3 ต่อ ไม่ต้องนับว่า 2 ต่อ