



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

What Is Pattern Discovery?

❑ What are patterns? *มส่ำน patterns ห้ข้อมูล / Data*

❑ **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**: *การค้นหาลักษณะเด่น* Uncovering patterns from massive data sets

❑ Motivation examples:

❑ What products were often purchased together? *สิ่ง ที่มักจะถูกซื้อด้วยกัน*

❑ What are the subsequent purchases after buying an iPad? *หลังจากซื้อ iPad แล้ว จะซื้ออะไรต่อ*

❑ What code segments likely contain copy-and-paste bugs?

❑ What word sequences likely form phrases in this corpus?

Pattern Discovery: Why Is It Important?

มีผลอย่างมาก เพราะเป็นพื้นฐานของ Data mining มากๆ

- ❑ 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

Basic Concepts: k-Itemsets and Their Supports

กลุ่ม items ที่รวมกันเข้าด้วยกัน

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

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

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

□ **(absolute) support (count)** of X, $\text{sup}\{X\}$:
Frequency or the number of occurrences of an itemset X

□ Ex. $\text{sup}\{\text{Beer}\} = 3$

□ Ex. $\text{sup}\{\text{Diaper}\} = 4$

□ Ex. $\text{sup}\{\text{Beer, Diaper}\} = 3$

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

→ การหา Beer คือการหา support ของ Beer

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. $s\{\text{Beer}\} = 3/5 = 60\%$

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

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

→ Beer มี 3 transaction = 60%

ดังนั้น support ของ Beer คือ 60% ของ transaction ทั้งหมด

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: $\frac{2}{5}$ (40%)

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. $Diaper \rightarrow Beer$ *คนที่ซื้อ Diaper มักจะซื้อ Beer ด้วย*

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 $X \cup Y$ *support คือ X และ Y*

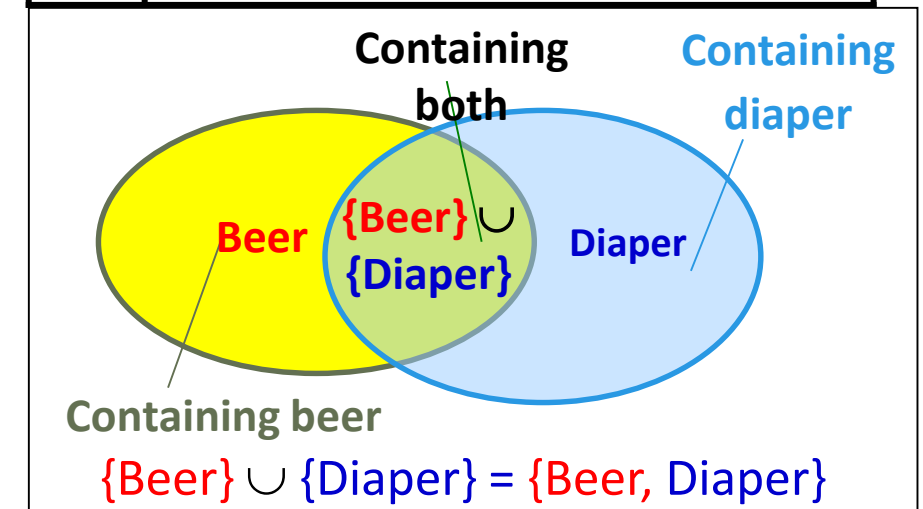
Ex. $s\{Diaper, Beer\} = 3/5 = 0.6$ (i.e., 60%)

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



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 $c \geq \text{minconf}$

□ Let $\text{minsup} = 50\%$

□ Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3

□ Freq. 2-itemsets: {Beer, Diaper}: 3

□ Let $\text{minconf} = 50\%$ $c = \frac{\text{sup}(X \cup Y)}{\text{sup}(X)}$

□ $\text{Beer} \rightarrow \text{Diaper}$ (60%, 100%)

□ $\text{Diaper} \rightarrow \text{Beer}$ (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

□ Observations:

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



ความสัมพันธ์



สูตร

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

Apriori Pruning and Scalable Mining Methods

หลักการตัดทิ้ง

หลักการ

- 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_k : Candidate itemset of size k

F_k : Frequent itemset of size k

$K := 1$;

$F_k := \{\text{frequent items}\}$; // frequent 1-itemset

While ($F_k \neq \emptyset$) **do** { // when F_k is non-empty

$C_{k+1} := \text{candidates generated from } F_k$; // candidate generation

Derive F_{k+1} by counting candidates in C_{k+1} with respect to TDB at minsup;

$k := k + 1$

}

return $\cup_k F_k$ // return F_k generated at each level

Pseudo-code = ขั้นตอนการหา frequent itemset ที่น่าสนใจ

The Apriori Algorithm—An Example

① none item set is removed in support less than 2

Database TDB

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

minsup = 2

C_1

1st scan

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

F_1

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

item set is removed in support less than 2

item set

C_2

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

2nd scan

C_2

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

F_2

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

item set

C_3

Itemset
{B, C, E}

3rd scan

F_3

Itemset	sup
{B, C, E}	2