



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

การทำ mining เพื่อหา Patterns ที่เกิดขึ้นซ้ำๆ

What Is Pattern Discovery?

- ❑ **What are patterns?** การค้นหา patterns ที่ซ่อนอยู่
ลูกค้านักซื้ออะไรคู่กันเสมอ (set of items)
- ❑ **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
- ❑ **Motivation examples:** เอาไปใช้ประโยชน์อะไรได้บ้าง?
 - ❑ What products were often purchased together? สินค้าใดที่ลูกค้ามักจะซื้อคู่กันเสมอ? เพราะถ้ารู้ว่า ร้านจะได้อะไรกับสินค้าที่ซื้อไว้เหมือนกัน
 - ❑ What are the subsequent purchases after buying an iPad? เมื่อลูกค้าซื้อ iPad ไปแล้ว จะกลับมาซื้อปากกา/เคส เป็นต้น
 - ❑ What code segments likely contain copy-and-paste bugs? ครั้งหนึ่งร้านอาจเจอโปรแกรมเมอร์ไม่สนใจลูกค้า
 - ❑ What word sequences likely form phrases in this corpus?

Pattern Discovery: Why Is It Important?

- ❑ 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

ตาราง, รูปภาพ, วิดีโอ, แผนที่

6

Basic Concepts: k-Itemsets and Their Supports

- ❑ **Itemset**: A set of one or more items
- ❑ **k-itemset**: $X = \{x_1, \dots, x_k\}$
 - ❑ 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. $\text{sup}\{\text{Beer, Eggs}\} = 1$

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

มักใช้ relative แทน absolute

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

จุดมุ่งหมายคือ เราต้องทำมือได้
เข้าใจการทำงานของ Data Mining

7

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 ดูยังไม่ว่ามีความถี่ = เกิดขึ้นบ่อย
- Let $\sigma = 50\%$ (σ : *minsup* threshold) ค่าขีดแบ่งว่า จะเอา หรือ ไม่เอา
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: coffee 2/5 (40%)
 - {Beer, Diaper}: 3/5 (60%)
 - All the frequent 3-itemsets?
 - None

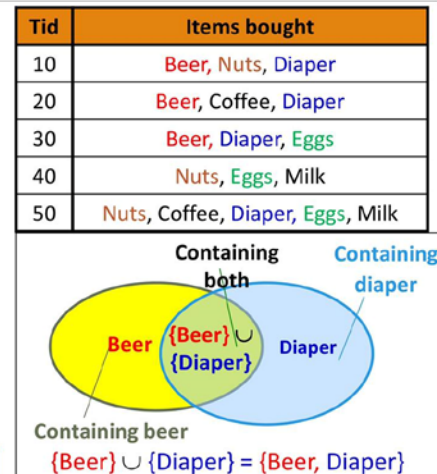
Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

- 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

8

From Frequent Itemsets to Association Rules

- Comparing with itemsets, rules can be more telling
 - Ex. *Diaper* \rightarrow *Beer* คนซื้อ Diaper จะนำไปสู่การซื้อ 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 $\frac{D/B}{D} = \frac{3/5}{4/5}$
 - Calculation: $c = \text{sup}(X \cup Y) / \text{sup}(X)$
 - Ex. $c = \text{sup}\{\text{Diaper}, \text{Beer}\} / \text{sup}\{\text{Diaper}\} = 3/4 = 0.75$



Note: $X \cup Y$: the union of two itemsets
 ■ The set contains both X and Y

9

Mining Frequent Itemsets and Association Rules

Association rule mining ก่อนจะหา ต้องกำหนด minsup, minconf

- Given two thresholds: **minsup, minconf**
- Find **all** of the rules, $X \rightarrow Y (s, c)$
 - such that, $s \geq \text{minsup}$ and $c \geq \text{minconf}$

อาจเลือกจากตัวที่มี support และ confident สูงสุด

Let **minsup = 50%**

- Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
- Freq. 2-itemsets: {Beer, Diaper}: 3

Let **minconf = 50%**

- $\text{Beer} \rightarrow \text{Diaper}$ (60%, 100%)
- $\text{Diaper} \rightarrow \text{Beer}$ (60%, 75%)

$$c = \frac{\text{sup}(X \cup Y)}{\text{sup}(X)}$$

(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

10

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

15

Apriori Pruning and Scalable Mining Methods

หลักการสำคัญ

- 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)

17

The Apriori Algorithm (Pseudo-Code)

C_k : Candidate itemset of size k

F_k : Frequent itemset of size k

```
K := 1;
Fk := {frequent items}; // frequent 1-itemset
While (Fk != ∅) do { // when Fk is non-empty
    Ck+1 := candidates generated from Fk; // candidate generation
    Derive Fk+1 by counting candidates in Ck+1 with respect to TDB at minsup;
    k := k + 1
}
return ∪k Fk // return Fk generated at each level
```

ข้อได้คือ = เป็นโค้ดโปรแกรม แต่ไม่ได้เป็นภาษาใดภาษาหนึ่ง เขียนในภาษานี้ไปแปลงเป็นภาษาที่เราใช้ได้

19

The Apriori Algorithm—An Example

