

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?

โดยทั่วไปแล้วจะหมายถึงสิ่งที่มักจะเกิดขึ้นร่วมกัน (patterns มักจะเกิดขึ้นซ้ำๆ), set of items มักจะเกิดขึ้นซ้ำๆ

❑ **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:

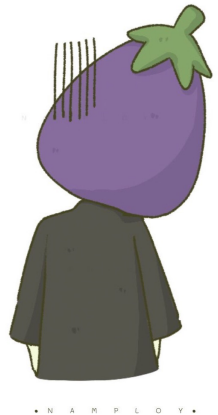
❑ What products were often purchased together? → สินค้าอะไรที่คนจะซื้อพร้อมกันเสมอ

❑ What are the subsequent purchases after buying an iPad? → หลังจากซื้อชิ้นนี้ไปแล้ว คราวหน้าจะมาซื้ออะไร

❑ What code segments likely contain copy-and-paste bugs? → มี code ของคนทำก่อนหน้ามาใช้ แก้ปัญหาได้ แต่จะมี bug เหมือนกัน

❑ What word sequences likely form phrases in this corpus?

↳ มันจะสามารถคาดเดาได้หรือไม่



Basic Concepts: k-Itemsets and Their Supports

งานที่มอบหมาย !!

- ☐ **Itemset**: A set of one or more items
เซตของไอเท็มที่คนมักซื้อพร้อมกัน
- ☐ **k-itemset**: $X = \{x_1, \dots, x_k\}$
ไอเท็มเซตที่มีสมาชิก k ตัว
- ☐ Ex. {Beer, Nuts, Diaper} is a 3-itemset
ไอเท็มเซตที่มี 3 รายการ
- ☐ **(absolute) support (count)** of X, $\text{sup}\{X\}$:
Frequency or the number of occurrences of an itemset X
จำนวนของ transaction ที่มา support X
- ☐ Ex. $\text{sup}\{\text{Beer}\} = 3$
จำนวน transaction ที่มี Beer
- ☐ Ex. $\text{sup}\{\text{Diaper}\} = 4$
จำนวน transaction ที่มี Diaper
- ☐ Ex. $\text{sup}\{\text{Beer, Diaper}\} = 3$
จำนวน transaction ที่มี Beer และ Diaper พร้อมกัน
- ☐ 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

งานที่มอบหมาย ID
คือ ทำงาน
ที่ 190
หรือ

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

- ☐ Ex. $s\{\text{Beer}\} = 3/5 = 60\%$
จากทั้งหมด 5 รายการ (จำนวน transaction ทั้งหมด)
- ☐ Ex. $s\{\text{Diaper}\} = 4/5 = 80\%$
- ☐ Ex. $s\{\text{Beer, Eggs}\} = 1/5 = 20\%$

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

ดูตามน้อย
ๆเกิด 50%.

- Let $\sigma = 50\%$ (σ : *minsup* threshold)
- For the given 5-transaction dataset

- All the frequent 1-itemsets:

Beer มี 3 ใน 5 transaction คือ 60%.

- Beer: 3/5 (60%); Nuts: 3/5 (60%)
- Diaper: 4/5 (80%); Eggs: 3/5 (60%)

- All the frequent 2-itemsets:

Beer กับ Diaper คู่กันอยู่ 3 ใน 5 transaction = 60%.

- {Beer, Diaper}: 3/5 (60%)
- All the frequent 3-itemsets?
- None

coffee : 2/5 (40%) => ไม่ผ่านเกณฑ์

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. $Diaper \rightarrow Beer$ → คนซื้อ Diaper มักไปซื้อ 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

รู้ความน่าจะเป็น support ของ X และ Y

Support, s: The probability that a transaction contains $X \cup Y$

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

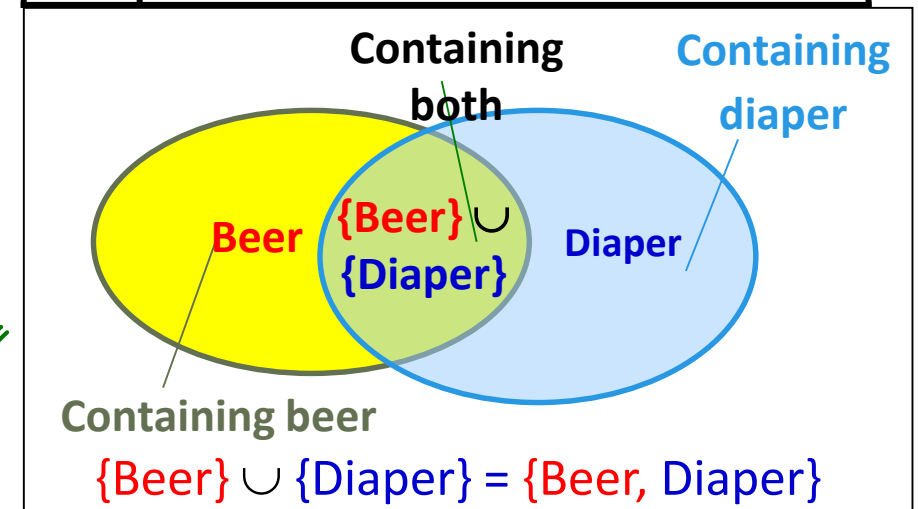
Confidence, c: The *conditional probability* that a transaction containing X also contains Y

หา support ของ X แล้วหารด้วย support ทั้งหมด

Calculation: $c = \text{sup}(X \cup Y) / \text{sup}(X)$

Ex. $c = \text{sup}\{Diaper, Beer\} / \text{sup}\{Diaper\} = \frac{3}{4} = 0.75$

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, $s \geq minsup$ and $c \geq minconf$

Let $minsup = 50\%$ ^{$minsupport$} \rightarrow เป็น 10% ที่มารีทที่น้อยเป็นส่วนใหญ่ของรายการแรกขึ้น

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

Let $minconf = 50\%$ ^{เกิดขึ้นซ้ำๆ กันมากแค่ไหน} $\frac{sup(Beer / Diaper)}{sup(Beer)}$

^{กฎ 2 ตัว} $\left\{ \begin{array}{l} Beer \rightarrow Diaper \text{ (60\%, 100\%)} \rightarrow \text{10% Diaper ไปวางขึ้นท้าย Beer} \\ Diaper \rightarrow Beer \text{ (60\%, 75\%)} \end{array} \right.$

(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

Efficient Pattern Mining Methods

- ❑ The Downward Closure Property of Frequent Patterns
- ❑ The ^{เอปรีออริ} ^{อัลกอริทึม} 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

ตัดแต่ง

ถ้าเกิดว่าไอเท็มเซตที่ค่าไม่ผ่าน MinSupport ไอเท็มเซตที่สูงกว่าก็ไม่มีทางผ่าน MinSupport (๑๙๙๔)

- ❑ 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

การเงื่อนไขผ่านขั้นตอน

- Outline of Apriori (level-wise, candidate generation and test)
 - Initially, scan DB once to get frequent 1-itemset → สแกนเริ่มต้นจากจุด 1 ไปที่ผลลัพธ์ก่อน
 - 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

มี 4 รายการ

กำหนด minsup
minsup = 2

C_1

1st สแกน
หา 1 ไอเท็มที่เข้า

มี 5 one Itemset

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

ไม่พอ minsup ไม่เอา

F_1

จะได้รายการใหม่ กำจัดที่ 1

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

กำจัด Itemset

C_2

หาทุกตัวที่ตรงกับ minsup หรือไม่

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

2nd scan

C_2

ได้

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

F_2

รายการใหม่ กำจัดที่ 2

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

จะได้

สร้าง tree Itemset

C_3

Itemset
{B, C, E}

จริงๆ ได้ 3 ตัว แต่เขาไม่นับเอาแล้ว

3rd scan

F_3

Itemset	sup
{B, C, E}	2

ทั้งหมดสุดท้าย

จบ ไปดูโปรแกรม
10.22 น.