

mean pattern ที่พบใน Data

# What Is Pattern Discovery?

## □ What are patterns?

รูปแบบบางอย่างที่เจอ

ข้อดีที่ง่าย

□ **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? สินค้าอะไรที่โดนมาก: ซื้อด้วยกันได้

□ What are the subsequent purchases after buying an iPad?

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

□ 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  
*in m n*
  - ❑ 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

ขั้นตอนการทำงาน

# Basic Concepts: k-Itemsets and Their Supports

สิ่งที่ต้องรู้ก่อน

คือ การหาว่า: ได้ สินค้าไหนบ้าง ตาม หลัก: ซึ่งรวมกัน

ที่ หาได้: ใน 1 ครั้ง (หรือ 1 ครั้ง/วัน) 2 ครั้ง 1 เดือน ฯลฯ

สินค้าที่ซื้อ  
คือ Itemset

**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,  $\text{sup}\{X\}$ :

จำนวนของรายการที่ซื้อสินค้าในชุด X

Frequency or the number of occurrences of an itemset X

Ex.  $\text{sup}\{\text{Beer}\} = 3$

สินค้าที่ซื้อได้ 3 รายการที่ซื้อ Beer

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

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

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

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

คือ จำนวน transactions ที่ซื้อสินค้าในชุด X หารด้วย จำนวน transactions ทั้งหมด

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

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

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

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

# 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  $\sigma$

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

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

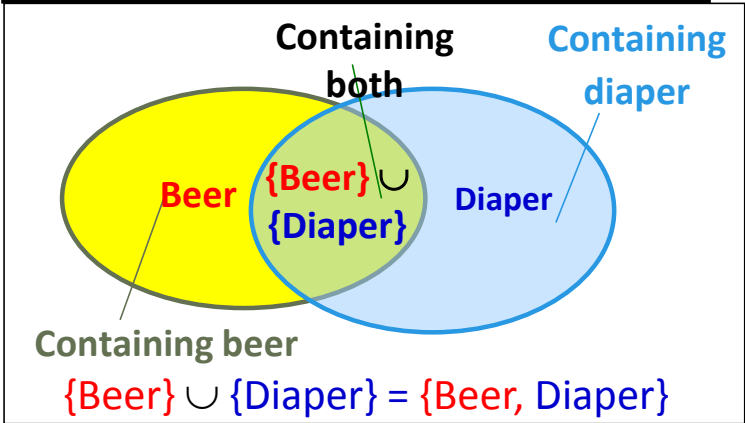
- All the frequent 1-itemsets:
  - Beer: 3/5 (60%); Nuts: 3/5 (60%)
  - Diaper: 4/5 (80%); Eggs: 3/5 (60%)
- All the frequent 2-itemsets:
  - {Beer, Diaper}: 3/5 (60%)
- All the frequent 3-itemsets?
  - None

- 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.  $Diaper \rightarrow Beer$  *ถ้าซื้อ Diaper ก็มักจะซื้อ 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  $X \cup Y$  *คือ X และ Y 出现在了กี่ครั้ง*
    - Ex.  $s\{Diaper, Beer\} = 3/5 = 0.6$  (i.e., 60%)
  - Confidence**,  $c$ : *for support of itemset* The **conditional probability** that a transaction containing  $X$  also contains  $Y$ 
    - Calculation:  $c = \frac{\text{sup}(X \cup Y)}{\text{sup}(X)}$
    - Ex.  $c = \frac{\text{sup}\{Diaper, Beer\}}{\text{sup}\{Diaper\}} = \frac{3}{4} = 0.75$

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,  $s \geq \text{minsup}$  and  $c \geq \text{minconf}$
- Let *minsup* = 50%
  - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
  - Freq. 2-itemsets: {Beer, Diaper}: 3
- Let *minconf* = 50%
  - *Beer*  $\rightarrow$  *Diaper* (60%, 100%)
  - *Diaper*  $\rightarrow$  *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

မိမိတို့ကလေးတွေ  
မိမိတို့ကလေးတွေ


# 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

# The Downward Closure Property of Frequent Patterns

---

- ❑ Observation: From TDB<sub>1</sub>: T<sub>1</sub>: {a<sub>1</sub>, ..., a<sub>50</sub>}; T<sub>2</sub>: {a<sub>1</sub>, ..., a<sub>100</sub>}
  - ❑ We get a frequent itemset: {a<sub>1</sub>, ..., a<sub>50</sub>}
  - ❑ Also, its subsets are all frequent: {a<sub>1</sub>}, {a<sub>2</sub>}, ..., {a<sub>50</sub>}, {a<sub>1</sub>, a<sub>2</sub>}, ..., {a<sub>1</sub>, ..., a<sub>49</sub>}, ...
  - ❑ There must be some hidden relationships among frequent patterns!
- ❑ The **downward closure (also called “Apriori”)** property of frequent patterns
  - ❑ If **{beer, diaper, nuts}** is frequent, so is **{beer, diaper}**
  - ❑ Every transaction containing {beer, diaper, nuts} also contains {beer, diaper}
  - ❑ Apriori: Any subset of a frequent itemset must be frequent
- ❑ Efficient mining methodology
  - ❑ If **any subset of an itemset S** is infrequent, then there is no chance for S to be frequent—why do we even have to consider S!?  A sharp knife for pruning!



# Apriori Pruning and Scalable Mining Methods

---

- Apriori pruning principle<sup>of course</sup>: 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 *အတူတူနဲ့*
    - Generate length-(k+1) candidate itemsets from length-k frequent itemsets *အတူ k+1 လောက် ဖန်တီး*
    - 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

เงื่อนไขที่สามารถนำไปประมวลผลได้

# The Apriori Algorithm—An Example

