# CS 412 Intro. to Data Mining

## Chapter 6. Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods
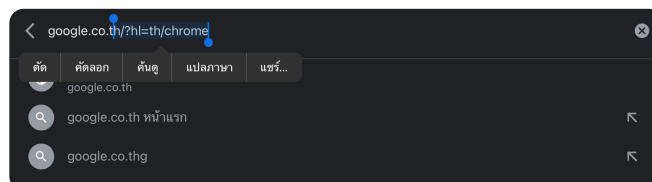
**Jiawei Han, Computer Science, Univ. Illinois at Urbana-Champaign, 2017**

*การทำ mining เพื่อหา Patterns ที่เกิดขึ้นซ้ำ*

---

# What Is Pattern Discovery?

- What are patterns? *การค้นหา patterns ที่ซ่อนอยู่*    *ลูกค้ามักซื้ออะไรคู่กันเสมอ*
    - **Patterns**: A set of items, subsequences, or substructures that occur *(set of items)* 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*    *สินค้าใดที่ลูกค้ามักจะซื้อคู่กันเสมอ ?*
    - What products were often purchased together? *เพราะ ลาภรัพย์ ร้านขายได้ เตรียมขายสินค้าที่ลูกจัดวางไว้เท่าๆกัน*
    - What are the subsequent purchases after buying an iPad? *เมื่อลูกค้า ซื้อ ไอแพด ไปแล้ว จากกล่อง ซื้อฝาฉก/เคส*
    - What code segments likely contain copy-and-paste bugs? *เงินต้น ตรังหน้าร้าน อาจจัดโปร โมชั่นให้ลูกค้า*
    - What word sequences likely form phrases in this corpus?

*ถ้ารวมข้อมูล ด้แนว แยะ ยังไม่เติบ ดิจ เช่น*



*จึน จะ ซิบ ดิท ให้ เค ต้อง การ ค้นทา ซึ่งเร็บ รู้จาก patterns เป็นต้น*

# 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

6

---

*สามารก เชื่อมเชิง*
*ต่อกันพ K ได้*

# Basic Concepts: k-Itemsets and Their Supports

❑ Itemset: A set of one or more items

❑ k-itemset: $X = \{x_1, ..., x_k\}$

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

❑ (absolute) support (count) of X, sup{X}: Frequency or the number of occurrences of an itemset X

*จำนวน transaction ที่มา สนับสนุน*

| 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*

   ❑ Ex. sup{Beer} = 3    *Absolute ไม่ดี เพราะไม่รู้*
   ❑ Ex. sup{Diaper} = 4   *transaction ทั้งหมด*
   ❑ Ex. sup{Beer, Diaper} = 3
   ❑ Ex. sup{Beer, Eggs} = 1

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

   ❑ Ex. s{Beer} = 3/5 = 60%
   ❑ Ex. s{Diaper} = 4/5 = 80%
   ❑ Ex. s{Beer, Eggs} = 1/5 = 20%

*จุดมุ่ง หมาย คือ เราต้อง ทำ ถึงได้*
*เข้าใจการทำงานของ Data mining*

# Basic Concepts: Frequent Itemsets (Patterns)

ดูยังไงว่าตัดตรงไหนดี - เกิดขึ้นบ่อย

- ❑ An itemset (or a pattern) X is *frequent* if the support of X is no less than a *minsup* threshold σ, ค่าขีดแบ่งว่า จะเอาหรือไม่เอา
- ❑ Let σ = *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: 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

---

# From Frequent Itemsets to Association Rules

- ❑ Comparing with itemsets, rules can be more telling
  - ❑ Ex. *Diaper → Beer* คนซื้อ Diaper จะน่าไปสู่การซื้อ 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
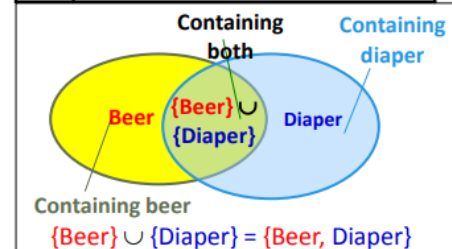  - ❑ Support, *s*: The probability that a transaction contains $X \cup Y$
  - ❑ Ex. s{Diaper, Beer} = 3/5 = 0.6 (i.e., 60%)
  - ❑ Confidence, *c: The conditional probability* that a transaction containing X also contains Y
  - ❑ Calculation: $c = sup(X \cup Y) / sup(X)$  D|B  D  3/5 ÷ 4/5
  - ❑ Ex. $c = sup\{Diaper, Beer\}/sup\{Diaper\} = ¾ = 0.75$

| 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 |

Containing both | Containing diaper

Beer   {Beer} ∪ {Diaper}   Diaper

Containing beer

{Beer} ∪ {Diaper} = {Beer, Diaper}

Note: X ∪ Y: the union of two itemsets
- The set contains both X and Y

# Mining Frequent Itemsets and Association Rules

*ผู้ใช้หมายต้องกำหนด minsup, minconf*

- ❑ **Association rule mining**
  - ❑ Given two thresholds: *minsup, minconf*
  - ❑ Find **all** of the rules, $X \rightarrow Y$ (s, c)
    - ❑ such that, $s \geq minsup$ and $c \geq minconf$

  *ตองเลือก จากตัวที่มี support และ confident สูงสุด*

- ❑ Let *minsup = 50%*
  - ❑ Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
  - ❑ Freq. 2-itemsets: {Beer, Diaper}: 3

  *สูตร*
  $$C = Sup(X \cup Y) / sup(X)$$

- ❑ Let *minconf = 50%*
  - ❑ *Beer → Diaper* (60%, 100%)
  - ❑ *Diaper → 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

# 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 $TDB_1$: $T_1$: $\{a_1, ..., a_{50}\}$; $T_2$: $\{a_1, ..., a_{100}\}$
  - ❑ We get a frequent itemset: $\{a_1, ..., a_{50}\}$
  - ❑ Also, its subsets are all frequent: $\{a_1\}, \{a_2\}, ..., \{a_{50}\}, \{a_1, a_2\}, ..., \{a_1, ..., a_{49}\}, ...$
  - ❑ 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

หลักการ ลำ ดับ

- ❑ 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—An Example



$m$ one itemset ว่ามีความเป็นไปได้อะไรบ้าง

$m$ support ของแต่ละ one itemset

∴ เนื่องจาก minsup = 2

เราจะคัด itemset ที่ sup ≤ 2

**Database TDB**

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

minsup = 2

$C_1$

1$^{st}$ 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 |

จะได้ตารางใหม่ดังนี้

ทุกนี้จะนับ support + ของแต่ละคู่

จับมาจับคู่กัน

แล้วนับ minsup มาตัดออกก็ได้ พวได้

$C_2$

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

2$^{nd}$ 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 |

ตั้งนั้น two itemset ทั้นะมาในรวมเครือ

$C_3$

| Itemset |
|---------|
| {B, C, E} |

3$^{rd}$ scan

$F_3$

| Itemset | sup |
|---------|-----|
| {B, C, E} | 2 |