# 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

---

# What Is Pattern Discovery?

- What are patterns? เป็นการค้นหา Patterns ที่ซ่อนอยู่

    - 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:

    สินค้าใดที่ลูกค้ามักจะซื้อคู่กันเสมอ ทำให้ร้านสามารถเตรียมของที่คู่กันไว้อย่างพอดี

    - What products were often purchased together?

    - 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?

5

# 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

---

# 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

   ❏ Ex. sup{Beer} = 3

   ❏ Ex. sup{Diaper} = 4

   ❏ Ex. sup{Beer, Diaper} = 3

   ❏ Ex. sup{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 |

❏ (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%

K-itemset ตัว k สามารถเปลี่ยนเป็นตัวเลขได้

Absolute support เป็นการนับจำนวน transaction ที่มาสนับสนุน แต่วิธีนี้ไม่รู้จำนวนทั้งหมดของข้อมูล

Relative support วิธีนี้เราจะสามารถรู้ถึงสัดส่วน และ จำนวนทั้งหมดของข้อมูลทั้งหมดด้วย

# Basic Concepts: Frequent Itemsets (Patterns)

ค่าความถี่ ดูว่าเหตุการณ์นี้เกิดขึ้นบ่อยแค่ไหน

- ❑ An itemset (or a pattern) X is *frequent* if the support of X is no less than a *minsup* threshold σ
- ❑ 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:
    - ❑ {Beer, Diaper}: 3/5 (60%)
  - ❑ All the frequent 3-itemsets?
    - ❑ None

Minsup threshold = ค่าขีดแบ่ง

| 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
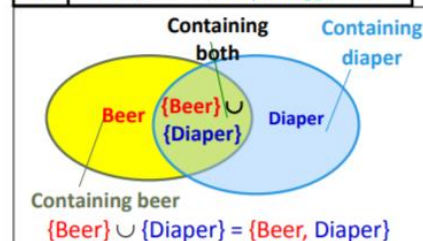
---

# From Frequent Itemsets to Association Rules

- ❑ Comparing with itemsets, rules can be more telling
  - ❑ Ex. *Diaper → Beer*   คนซื้อ Diaper จะนำไปสู่การ ซื้อ 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: $c = \sup(X \cup Y) / \sup(X)$
  - ❑ Ex. $c = \sup\{Diaper, Beer\}/\sup\{Diaper\} = ¾ = 0.75$

(D,B) / (D)

(3/5) /(4/5)

| 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

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

- Let *minconf = 50%*
  - *Beer → Diaper* (60%, 100%)
  - *Diaper → Beer* (60%, 75%)

  (Q: Are these all rules?)

*ต้องมีการกำหนด minsup, minconf*

| 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

- ❑ 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

ซูโดโค้ด เป็นโค้ดโปรแกรม แต่ไม่ได้ภาษาใดภาษาหนึ่ง
เขียนขึ้นเพื่อให้เรานำไปแปลงไปเป็นภาษาที่เราใช้ได้

K := 1;

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

**While** ($F_k$ != ∅) **do** {     // when $F_k$ is non-empty

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

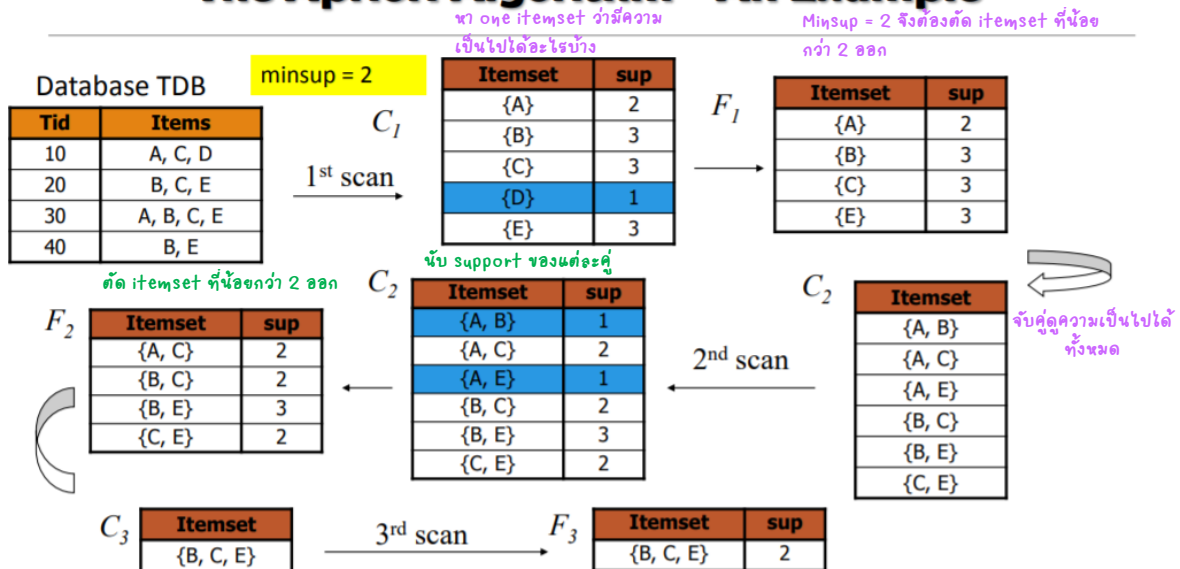   Derive $F_{k+1}$ by counting candidates in $C_{k+1}$ with respect to *TDB* at minsup;

   k := k + 1

   **}**

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

# The Apriori Algorithm—An Example