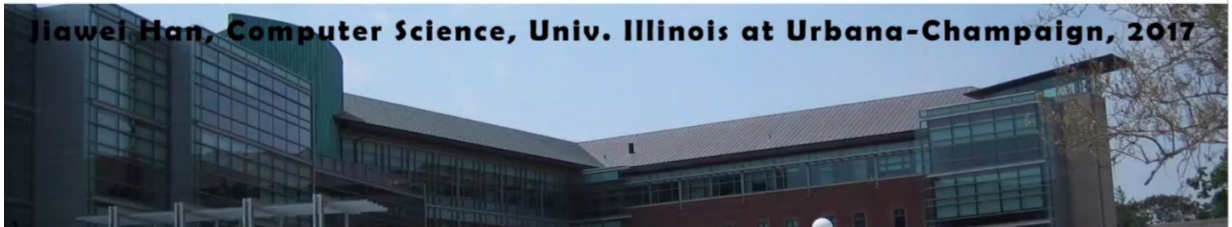# CS 412 Intro. to Data Mining

## 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?**

  ❑ **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 *การค้นหารูปแบบ*

❑ Motivation examples: *Ex. เราซื้อไอแพดไปแล้ว ครั้งถัดไปเขาจะซื้ออะไร ↳ เริ่มจากก*

  ❑ 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?

---

# 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

เซ็ตของไอเท็ม : สินค้าในบิล ที่คนมักจะซื้อร่วมกัน

- **Itemset**: A set of one or more items
  จำนวนเติมตัวเลขอะไรลงไปก็ได้
- **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  เราไม่สามารถรู้ ว่ามันมากหรือน้อย
  - Ex. sup{Diaper} = ④  จากใน transaction
  - 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%  ❀ มีการบอกสัดส่วน ขอ 100 transaction ทั้งหมด
  - Ex. s{Diaper} = 4/5 = 80%
  - Ex. s{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
- Let σ = 50% (σ: *minsup* threshold)  ค่าตัดแบ่ง
  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%)
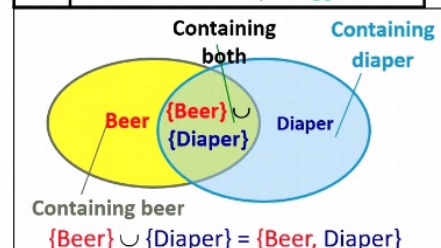  - All the frequent 3-itemsets?
  - None
  ถ้ามีค่าดีเกินแสดงว่ามีความน่าสนใจ

| 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*
    - *Buying diapers may likely lead to buying beers*
- How strong is this rule? (support, confidence)
  - Measuring association rules: $X \to Y$ (s, c)
  - Both X and Y are itemsets  X นำไปสู่การซื้อ Y
  - **Support, s:** The probability that a transaction contains $X \cup Y$  X และ Y ถูกซื้อพร้อมกัน
    - 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. $c = sup\{Diaper, Beer\}/sup\{Diaper\} = \frac{3}{4} = 0.75$

| 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 |

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

คือ Confidence ของ Rule การซื้อของ

# 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%*    $\tilde{\text{3}}\text{ठ}\text{ज़}$ $C = \frac{sup (X \cup Y)}{sup (x)}$
  ❑ *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 Apriori Algorithm (Pseudo-Code)

$C_k$: Candidate itemset of size k

$F_k$ : Frequent itemset of size k

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

# 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

# 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)

# The Apriori Algorithm—An Example

Database TDB

minsup = 2

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

1st scan

$C_1$

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

$F_1$

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

$C_2$

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

$F_2$

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

2nd scan

$C_2$

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

$C_3$

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

3rd scan

$F_3$

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