# 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
- Motivation examples:
  - 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

- Itemset: A (set) of one or more items    *ชุดที่ซื้อร่วมกัน (รถ)*
- k-itemset:  X = {x1, …, xk}    *ซื้อมาใน 3 ตัว*
  - 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    *มีใน transactionนี้มี beerกี่*
  - Ex. sup{Diaper} = 4
  - Ex. sup{Beer, Diaper} = 3    *2 item set*
  - Ex. sup{Beer, Eggs} = 1

*√ คำนวณรอบ transactions ที่ สนับสนุน ความคิดตา*

| Tid | Items bought |
|-----|--------------|
| 10  | *transaction 1* Beer, Nuts, Diaper |
| 20  | *" 2* Beer, Coffee, Diaper |
| 30  | *" 3* Beer, Diaper, Eggs |
| 40  | *" 4* Nuts, Eggs, Milk |
| 50  | *" 5* 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)    *จาก transaction ทั้งหมดมีกี่ ใน 1 ชุด )*
  - Ex.  s{Beer} = 3/5 = 60%
  - Ex.  s{Diaper} = 4/5 = 80%
  - Ex.  s{Beer, Eggs} = 1/5 = 20%

---

*(threshold) ต้องมากกว่า 'ไว' /ไม่ไว , เอา Min*

# Basic Concepts: Frequent Itemsets (Patterns)

- An itemset (or a pattern) X is *frequent*    *บ่อย*
  if the support of X is no less than a *minsup* threshold σ    *เกิน 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

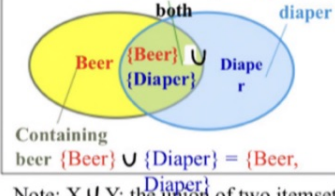| Tid | Items bought |
|-----|--------------|
| 10  | Beer, Nuts, Diaper |
| 20  | Beer, Coffee, Diaper |
| 30  | Beer, Diaper, Eggs |
| 40  | Nuts, Eggs, Milk |
| 50  | Nuts, Coffee, Diaper, Eggs, Milk |

- 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

*Coffee : 2/5 (40%) ไม่ถึง thresshold*

# From Frequent Itemsets to Association Rules

- Comparing with itemsets, rules can be more telling
  - Ex. *Diaper → Beer* ถ้า Diaper จะถึง Beer ด้วย
    - *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
  - **Support**, s: The probability that a transaction contains $X \cup Y$ → ร้อยละอะไร
    - 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. $c = \sup\{Diaper, Beer\}/\sup\{Diaper\} = \frac{3}{4} =$

| 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** พวกมาน= itemset ที่ใช้ใน
  - Given two thresholds: *minsup, minconf*
  - Find **all** of the rules, $X \to Y$ (s, c)
  - such that, s ≥ *minsup* and c ≥ *minconf*
- Let *minsup = 50%* minsup ตั้งไว้ 50%
  - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
  - Freq. 2-itemsets: {Beer, Diaper}: 3
  - → item เซ็ทที่อยู่บ่อยในฐานข้อมูล ของ itemset ที่ผ่านมา
- Let *minconf = 50%* เทิกว่านาน แต่ไม่
  - *Beer → Diaper* (60%, 100%) → Sup(Beer ∪ Diaper)
  - *Diaper → Beer* (60%, 75%) → Sup(Beer)

| 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

# Chapter 6: Mining Frequent Patterns, Association and Correlations: Basic Concepts and Methods

- Basic Concepts

- Efficient Pattern Mining Methods  ✗

- Pattern Evaluation

- Summary

# 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

# Apriori Pruning and Scalable Mining Methods

ตัวแปรง
- 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)

ถ้าเกิด item set
ที่เล็กกว่าไม่ชน
จะตัดทิ้ง

# Apriori: A Candidate Generation & Test Approach

- Outline of Apriori (level-wise, candidate generation and test) → *เขียนเป็นขั้นตอน*
  - Initially, scan DB once to get frequent 1-itemset
- Repeat     *สร้างได้เท่าไหร่*
  - 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   minsup = 2

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

1st scan

C1

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

F1

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

C2

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

2nd scan

C2

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

F2

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

C3

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

3rd scan

F3

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