

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?

คือไอเทมหลายๆอย่างรวมกันถูกซื้อพร้อมกัน เพื่อเอาไปใช้ประโยชน์ในเรื่องสินค้าตัวไหนที่ลูกค้ามักซื้อพร้อมกัน

Basic Concepts: k-Itemsets and Their Supports

- 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 , $\text{sup}(X)$: Frequency or the number of occurrences of an itemset X
 - Ex. $\text{sup}(\text{Beer}) = 3$
 - Ex. $\text{sup}(\text{Diaper}) = 4$
 - Ex. $\text{sup}(\text{Beer, Diaper}) = 3$
 - Ex. $\text{sup}(\text{Beer, Eggs}) = 1$

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. $s(\text{Beer}) = 3/5 = 60\%$
 - Ex. $s(\text{Diaper}) = 4/5 = 80\%$
 - Ex. $s(\text{Beer, Eggs}) = 1/5 = 20\%$

Basic Concepts: Frequent Itemsets (Patterns)

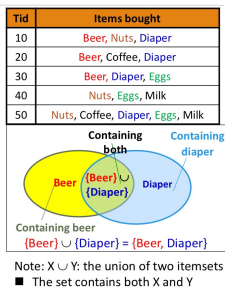
- An itemset (or a pattern) X is **frequent** if the support of X is no less than a **minsup** threshold σ
- Let $\sigma = 50\%$ (σ : minsup threshold) For the given 5-transaction dataset
 - All the frequent 1-itemsets:
 - Beer: 3/5 (60%); Nuts: 3/5 (60%)
 - Diaper: 4/5 (80%); Eggs: 3/5 (60%)
 - All 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

From Frequent Itemsets to Association Rules

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



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%
 - $\text{Beer} \rightarrow \text{Diaper}$ (60%, 100%)
 - $\text{Diaper} \rightarrow \text{Beer}$ (60%, 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

- Observations:**
 - Mining association rules and mining frequent patterns are very close problems
 - Scalable methods are needed for mining large datasets

การสกัดหารูปแบบซ้ำๆ ในข้อมูล โดยจากข้อมูลด้านบนสามารถสรุปได้ว่า คนที่ซื้อผ้าอ้อมกับเบียร์ มักจะเป็นคุณพ่อ เป็นแพทเทิร์นที่เกิดขึ้นซ้ำๆ เพื่อนำไปใช้ในการทำธุรกิจ