

# Basic Data Mining Algorithms

Liyao Xiang

<http://xiangliyao.cn/>

Shanghai Jiao Tong University

<http://jhc.sjtu.edu.cn/public/courses/EE226/>

# Notice

- There will be a quiz in the next week's class. Please take a piece of paper and pens.

# Reference and Acknowledgement

- Most of the slides are credited to Prof. Jiawei Han's book "Data Mining: Concepts and Techniques."

# Outline

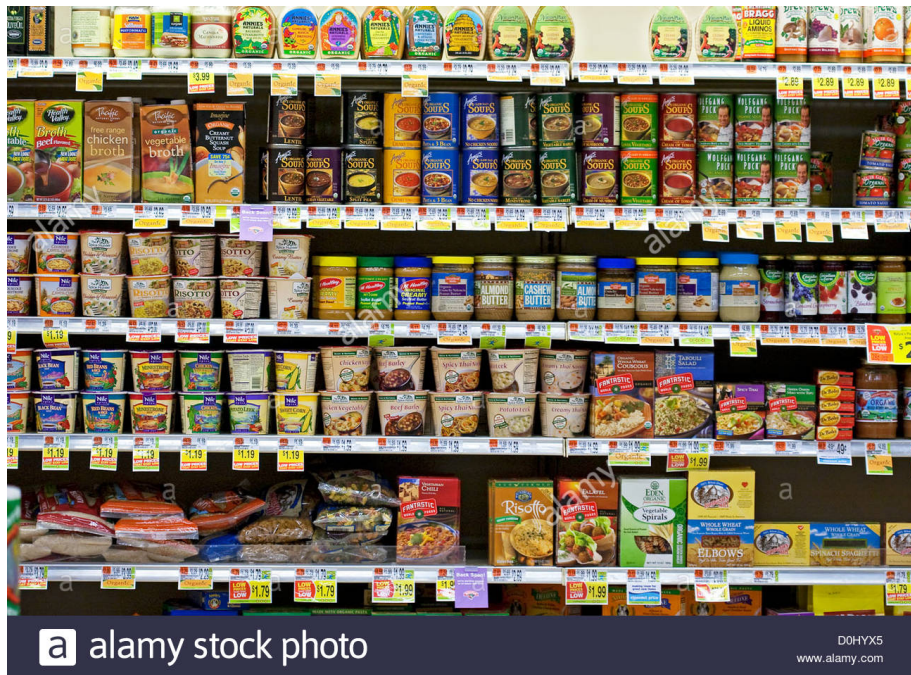
- Basic Concepts in Frequent Pattern Mining
- Frequent Itemset Mining Methods
- Pattern Evaluation Methods

# Outline

- Basic Concepts in Frequent Pattern Mining
- Frequent Itemset Mining Methods
- Pattern Evaluation Methods

# Basic Concepts

- **Frequent pattern:** a pattern (a set of items, subsequences, substructures ...) that appear frequently in a database

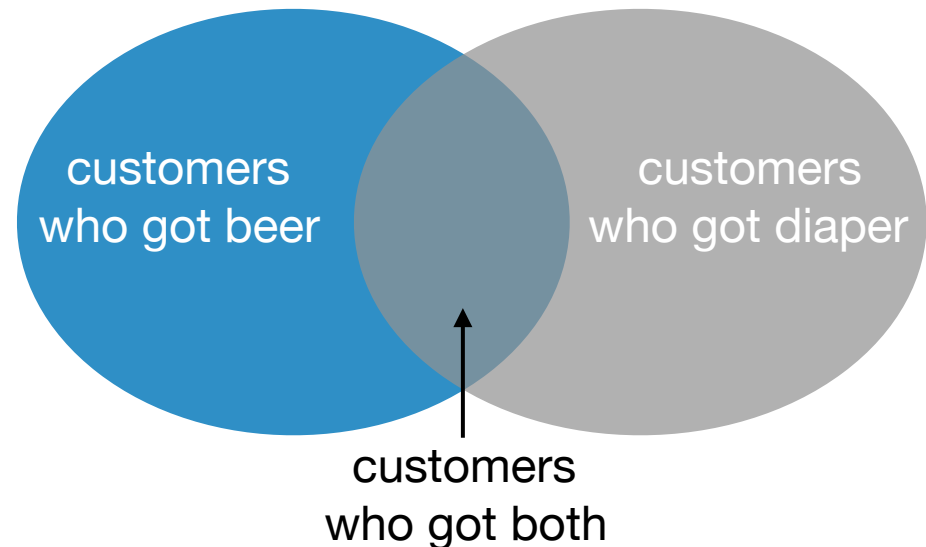


- Finding frequent patterns is key to mining associations, correlations, clustering, classification and other relationships among data.
- Applications: basket data analysis, cross-marketing, catalog design ...

# Basic Concepts

- **itemset**: a set of one or more items
- **k-itemset**:  $X = \{x_1, \dots, x_k\}$
- **(absolute) support**, or **support count** of  $X$ : frequency or occurrence of an itemset  $X$
- **(relative) support**: the fraction of transactions that contains  $X$  over all transaction
- An itemset  $X$  is **frequent** if  $X$ 's support is no less than a defined threshold **min\_sup**

TID	Items Purchased
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk



# Basic Concepts

- **support**: probability that a transaction contains  $X \cup Y$

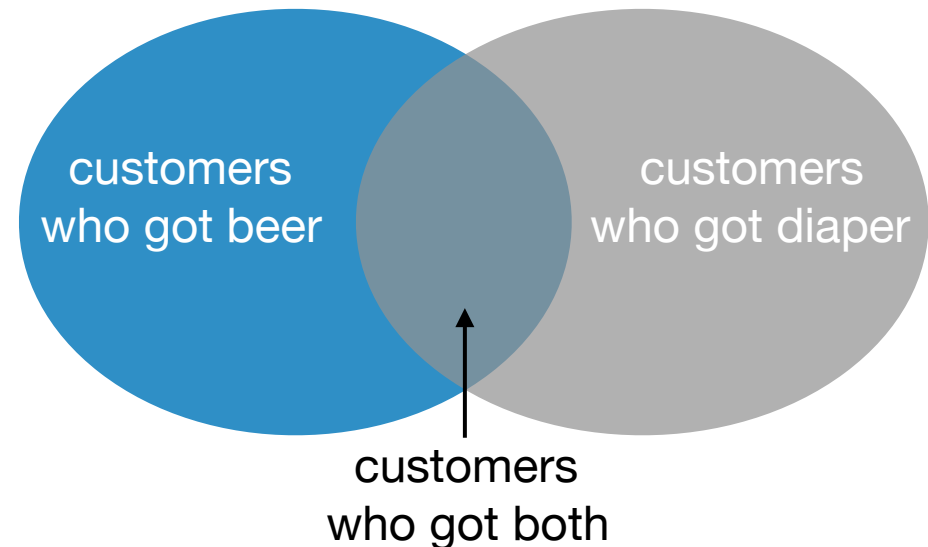
$$\text{support}(X \Rightarrow Y) = P(X \cup Y)$$

- **confidence**: conditional prob. that a transaction having  $X$  also contains  $Y$

$$\text{confidence}(X \Rightarrow Y) = P(Y|X)$$

$$P(Y|X) = \frac{\text{support}(X \cup Y)}{\text{support}(X)}$$

TID	Items Purchased
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk





# Basic Concepts

- **min\_sup**: minimum support threshold
- **min\_conf**: minimum support confidence threshold
- e.g., find all rules  $X \Rightarrow Y$  with min\_sup and min\_conf

let min\_sup = 50%, min\_conf = 50%

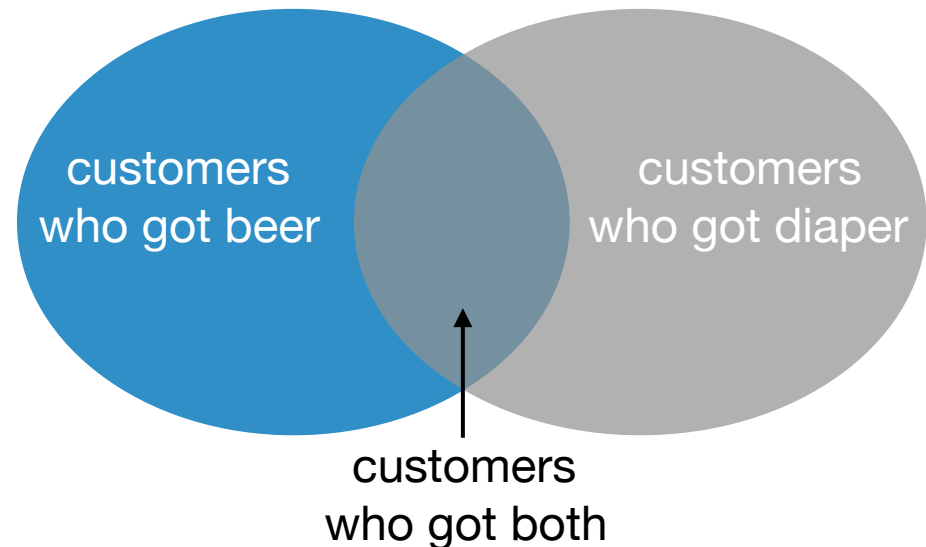
frequent pattern: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3, {Beer, Diaper}: 3

- Association rules:

Beer  $\Rightarrow$  Diaper (60%, 100%)

Diaper  $\Rightarrow$  Beer (60%, 75%)

TID	Items Purchased
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk



# Basic Concepts

- Association rule mining includes:
  1. **Find all frequent itemsets**: frequency of itemsets  $\geq \text{min\_sup}$
  2. **Generate strong association rules** from the frequent itemsets
- 1 is the major step, but challenging in that there may be **a huge number of** itemsets satisfying  $\text{min\_sup}$
- An itemset is frequent  $\Rightarrow$  each of its subsets is frequent
- Solution: mine **closed frequent itemset** and **maximal frequent itemset**
- **closed frequent itemset** X: X is frequent and there is no super-itemset  $Y \supset X$  with the same support count as X
  - closed frequent itemset is a lossless compression of frequent itemset
- **maximal frequent itemset** X: X is frequent and there is no super-itemset  $Y \supset X$  which is frequent

# Basic Concepts

- e.g.,  $\{\langle a_1, \dots, a_{100} \rangle, \langle a_1, \dots, a_{50} \rangle\}$ ,  $\text{min\_sup} = 1$
- What is the set of closed frequent itemset?
  - $\langle a_1, \dots, a_{100} \rangle$ : 1,  $\langle a_1, \dots, a_{50} \rangle$ : 2
- What is the set of maximal frequent itemset?
  - $\langle a_1, \dots, a_{100} \rangle$ : 1
- We can assert  $\langle a_2, a_{45} \rangle$  is frequent since  $a_2, a_{45} \in \langle a_1, \dots, a_{50} \rangle$  but cannot assert their actual support count
- How many itemsets are potentially to be generated in the worst case?
  - When  $\text{min\_sup}$  is low, there exist potentially an exponential number of frequent itemsets
  - Worst case:  $M^N$  where  $M = \#$  distinct items,  $N = \text{max length of transactions}$

# Summary

- frequent pattern
- k-itemset
- (absolute) support, support count, relative support
- min\_sup, confidence
- closed frequent itemset, maximal frequent itemset

# Outline

- Basic Concepts in Frequent Pattern Mining
- Frequent Itemset Mining Methods
- Pattern Evaluation Methods

# Frequent Itemset Mining Methods

- Apriori: A Candidate Generation-and-Test Approach
- Improving the Efficiency of Apriori
- FP-Growth: A Frequent Pattern-Growth Approach
- ECLAT: Frequent Pattern Mining with Vertical Data Format

# Apriori

- **Downward Closure** Property: any subset of a frequent itemset must be frequent
    - e.g., if {beer, diaper, nuts} is frequent, so is {beer, diaper} since every transaction having {beer, diaper, nuts} also contains {beer, diaper}
  - **Apriori** employs a level-wise search where k-itemsets are used to explore (k + 1)-itemsets. Steps:
    1. Scan database once to get frequent 1-itemsets  $L_1$
    2. **Join** the k-frequent itemsets  $L_k$  to **generate** length (k+1) candidate itemsets  $C'_{k+1}$
    3. **Prune**  $C'_{k+1}$  against the database to get  $C_{k+1}$
    4. Scan (**Test**) database for the count of each candidate in  $C_{k+1}$ , obtain  $L_{k+1}$
    5. Terminate when no frequent or candidate set can be generated
- 