

Chapter 3

Association Mining

Associate Professor Yachai Limpiyakorn, Ph.D.

Basic Concepts

TID	Produce
1	MILK, BREAD, EGGS
2	BREAD, SUGAR
3	BREAD, CEREAL
4	MILK, BREAD, SUGAR
5	MILK, CEREAL
6	BREAD, CEREAL
7	MILK, CEREAL
8	MILK, BREAD, CEREAL, EGGS
9	MILK, BREAD, CEREAL

- Given:
 - (1) database of transactions/ transactional database
 - (2) each transaction is a list of items purchased

- Find:

ความสัมพันธ์ที่น่าสนใจระหว่างไอเทมเซต (itemset) ในชุดข้อมูล ความสัมพันธ์ที่ได้เขียนอยู่ในรูปกฎความสัมพันธ์ (Association Rule) ของเซตของไอเทมที่เป็นเหตุ (Antecedent) ไปสู่เซตของไอเทมที่เป็นผล (Consequent)

{Cheese, Milk} → Bread [S=5%, C=80%]

80% of customers who buy cheese and milk also buy bread and 5% of customers buy all these products together



How can association rules be used?

Stories – Beer and Diapers

- ♦ **Diapers and Beer.** Most famous example of market basket analysis for the last few years. If you buy diapers, you tend to buy beer.
- T. Blischok headed Terradata's Industry Consulting group.
- K. Heath ran self joins in SQL (1990), trying to find two itemsets that have baby items, which are particularly profitable.
- Found this pattern in their data of 50 stores/90 day period.
- Unlikely to be significant, but it's a nice example that explains associations well.

Ronny Kohavi ICML 1998



Probably mom was calling dad at work to buy diapers on way home and he decided to buy a six-pack as well.

The retailer could move diapers and beers to separate places and position high-profit items of interest to young fathers along the path.

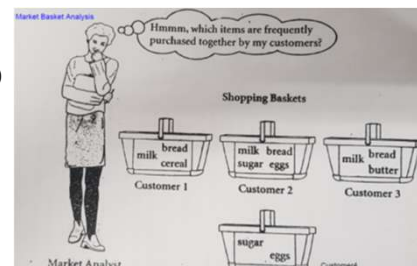
2110773-3 2/67

3

Application ₁

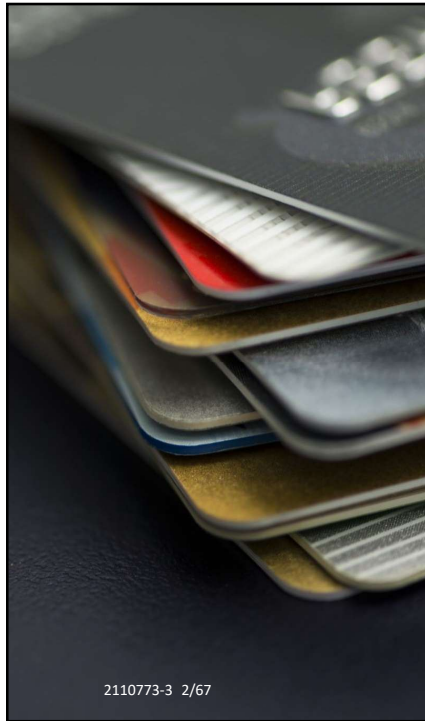
ส่วนใหญ่มักประยุกต์ใช้เทคนิคการทำเหมืองความสัมพันธ์กับการวิเคราะห์ทางการตลาด (Market Basket Analysis: MBA) ซึ่งเป็นรูปแบบการจัดกลุ่ม (Clustering) แบบหนึ่ง ที่ใช้เพื่อหากลุ่มสิ่งของที่น่าจะปรากฏร่วมกันในทรานแซกชันหนึ่งๆ มักเป็นทรานแซกชัน ณ จุดขาย (point-of-sale) ผลลัพธ์หรือแบบจำลองที่ได้สามารถแสดงได้ด้วยกฎซึ่งบอกความเป็นไปได้ของการซื้อผลิตภัณฑ์ต่างๆร่วมกัน การวิเคราะห์ทางการตลาดมีบทบาทสำคัญต่ออุตสาหกรรมการค้าปลีก (Retail industry) เพื่อให้ทราบถึงพฤติกรรมการซื้อสินค้าของลูกค้าซึ่งเป็นประโยชน์ในการ

- ♦ จัดพื้นที่ร้านค้า (Store layout)
- ♦ ทำตลาดเพื่อส่งเสริมการขายสินค้าหรือบริการซึ่งกันและกัน (Cross-marketing)
- ♦ ออกแบบหนังสือแคตตาล็อกสินค้า (Catalog design)
- ♦ วางแผนการส่งเสริมการขายและการตั้งราคาผลิตภัณฑ์ (Product pricing and promotion)



2110773-3 2/67

4



Application₂

นอกจากนี้ สามารถประยุกต์ใช้การวิเคราะห์ทางการตลาดกับกิจกรรมใกล้เคียงที่ลูกค้ามักกระทำด้วยกัน เพื่อก่อให้เกิดรายได้สูงสุดจากการจัดประเภทผลิตภัณฑ์หรือบริการเข้าด้วยกัน ได้แก่

- ◆ การใช้จ่ายผ่านบัตรเครดิตของลูกค้าในการเข้าพักโรงแรม เช่ารถ ทำให้สามารถทำนายค่าใช้จ่ายต่อไปของลูกค้า
- ◆ แพ็กเกจการให้บริการการสื่อสารโทรคมนาคม เพื่อก่อให้เกิดรายได้สูงสุด
- ◆ การให้บริการทางธนาคารที่ลูกค้ามักซื้อด้วยกัน เพื่อก่อให้เกิดประโยชน์สูงสุด เช่น ประเภอบัญชีที่ลูกค้ามักเปิดด้วยกัน (account bundle) การให้บริการการลงทุนครบวงจร และแพ็กเกจสินเชื่อการซื้อรถ เป็นต้น

2110773-3 2/67

5

นิยามพื้นฐาน การทำเหมือง ความสัมพันธ์

- ◆ ไอเทมเซต (itemset - I) คือเซตที่มีไอเทมทั้งหมดเป็นสมาชิก ซึ่งไอเทมในที่นี้อาจเป็นชื่อสินค้า หรือชื่อใดๆ ที่เป็นหน่วยพื้นฐานที่จะนำมาทำการเรียนรู้
- ◆ ทรานแซกชัน (transaction - T) เป็นเซตของไอเทม โดยที่ $T \subseteq I$
- ◆ เซตข้อมูล (data set - D) คือเซตที่มีทรานแซกชันทุกตัวเป็นสมาชิก
เรากล่าวว่าทรานแซกชัน T บรรจุเซตย่อยของไอเทม X ก็ต่อเมื่อ $X \subseteq T$
เพราะฉะนั้นจึงนิยามกฎความสัมพันธ์ได้ว่า
- ◆ กฎความสัมพันธ์ (Association Rule) คือการอุปนัยในรูปแบบ $X \rightarrow Y$ เมื่อ $X \subset I$, $Y \subset I$ และ $X \cap Y = \emptyset$

2110773-3 2/67

6

Objective measures of rule interest

- Support
- Confidence or strength
- Lift or Interest or Correlation
- Conviction
- Leverage or Piatetsky-Shapiro
- Coverage

Association Rule

- Rule form

Antecedent \rightarrow Consequent [*support*, *confidence*]

Note: *support* and *confidence* are user defined measures of interestingness

- Examples

$\text{buys}(x, \text{"computer"}) \rightarrow \text{buys}(x, \text{"financial management software"}) [0.5\%, 60\%]$

$\text{age}(x, \text{"30..39"}) \wedge \text{income}(x, \text{"42..48K"}) \rightarrow \text{buys}(x, \text{"car"}) [1\%, 75\%]$

Rule basic Measures: Support and Confidence

$$A \Rightarrow B [s, c]$$

Support: denotes the frequency of the rule within transactions. A high value means that the rule involve a great part of database.

$$\text{support}(A \Rightarrow B) = p(A \cup B)$$

Confidence: denotes the percentage of transactions containing A which contain also B. It is an estimation of conditioned probability .

$$\text{confidence}(A \Rightarrow B [s, c]) = p(B|A) = \text{sup}(A,B)/\text{sup}(A)$$

Calculation of Support and Confidence

• Support

คำนวณค่าสนับสนุน ได้จากจำนวนทรานแซกชันที่มีรายการ X และ Y เกิดร่วมกันหารด้วยจำนวนทรานแซกชันทั้งหมด

$$\begin{aligned}\text{support}(X \rightarrow Y) \\ &= P(X \cup Y) \\ &= \text{tran_count}(X \cup Y) / \text{tran_count}(D)\end{aligned}$$

• Confidence

คำนวณค่าความเชื่อมั่นได้จากจำนวน ทรานแซกชันที่มีรายการ X และ Y เกิดร่วมกันหารด้วยจำนวนทรานแซกชันที่มีรายการ X

$$\begin{aligned}\text{confidence}(X \rightarrow Y) \\ &= P(Y|X) \\ &= \text{tran_count}(X \cup Y) / \text{tran_count}(X)\end{aligned}$$

Practice Calculating Support and Confidence

Transaction ID	Items Bought
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F

ก. ให้คำนวณหาค่า support และ confidence ของความสัมพันธ์ $A \rightarrow C$ และ $C \rightarrow A$

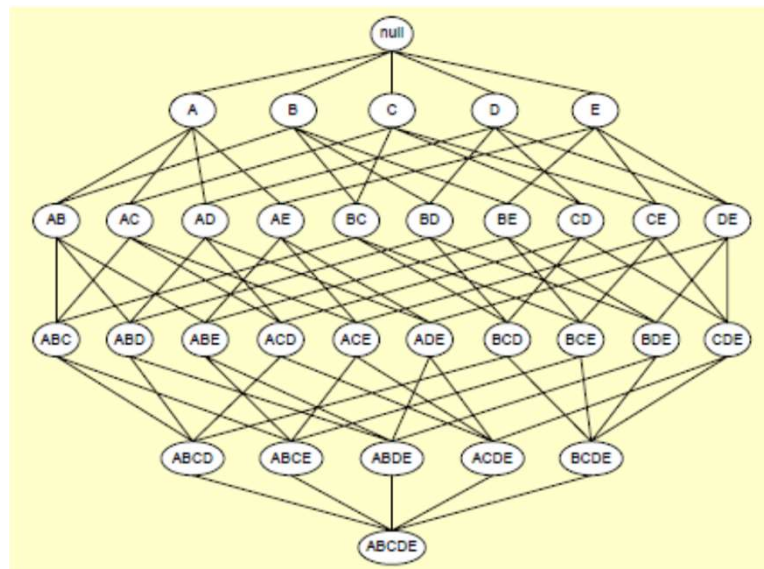
ข. กำหนดให้ minimum support = 50% และ minimum confidence = 80% อยากทราบว่าความสัมพันธ์ $A \rightarrow C$ และ $C \rightarrow A$ ความสัมพันธ์ใดเป็นกฎความสัมพันธ์

Association Mining

เป็นปัญหาการค้นหากฎความสัมพันธ์ นิยามได้ดังนี้

- **การค้นหากฎความสัมพันธ์** คือ การหาความสัมพันธ์ทั้งหมดในทรานแซกชันทุกตัวของเซตข้อมูลที่กำหนดให้ โดยกฎความสัมพันธ์ที่หาได้ทั้งหมดจะต้องมีค่าสนับสนุน (support) ไม่ต่ำกว่าค่าสนับสนุนน้อยสุด (minimum support) ที่ผู้ใช้กำหนดไว้ และมีค่าความเชื่อมั่น (confidence) ไม่ต่ำกว่าค่าความเชื่อมั่นน้อยสุด (minimum confidence) ที่ผู้ใช้ได้กำหนดไว้
- การค้นหากฎความสัมพันธ์สามารถแบ่งย่อยได้เป็นสองขั้นตอน คือ
 1. ค้นหาเซตของไอเทมปรากฏบ่อย (frequent itemset) หรือไอเทมเซตที่มีค่าสนับสนุนไม่ต่ำกว่าค่าสนับสนุนน้อยสุดที่กำหนดให้
 2. นำไอเทมเซตปรากฏบ่อยเหล่านั้นมาสร้างเป็นกฎความสัมพันธ์ต่อไป

Itemset Lattice



2110773-3 2/67

13

Apriori Principle

Any subset of a frequent itemset must also be frequent

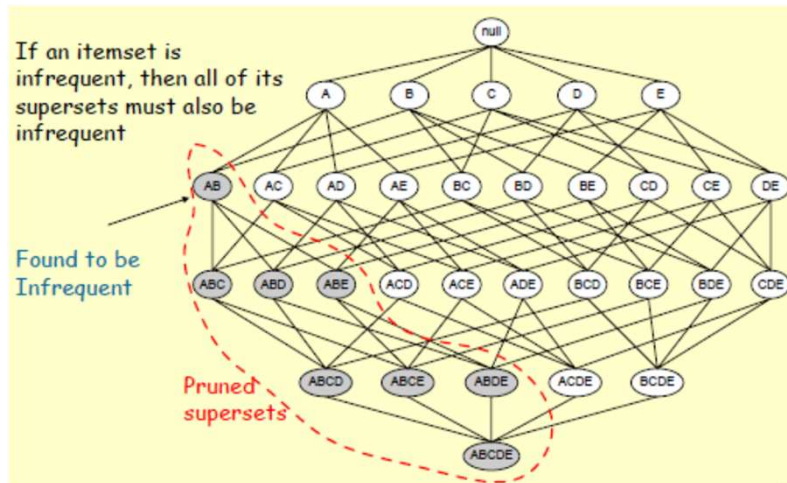
No superset of any infrequent itemset should be generated or tested

- Many item combinations can be pruned

2110773-3 2/67

14

Apriori Principle for Pruning Candidates



2110773-3 2/67

15

Association Mining: 2 key steps

1. Find all Frequent Itemsets: the sets of items that pass minimum support
 - ❖ Apriori Algorithm
 - มีการจัดเรียงลำดับของไอเทมในแต่ละทรานแซกชันก่อนประมวลผล
 - การสร้างไอเทมเซตจะสร้างตามระดับชั้น จากชั้นที่ k , $k+1$, $k+2$
 - ใช้ความรู้ก่อนหน้าคือคุณสมบัติของไอเทมเซตเกิดบ่อยในการตัดเล็ม
2. For every frequent itemset X , generate all non-empty subset S of X

$$S \rightarrow (X-S)$$

Output the rule $S \rightarrow (X-S)$
If confidence \geq min_confidence

2110773-3 2/67

16

Apriori Algorithm

Algorithm: Apriori. Find frequent itemsets using an iterative level-wise approach based on candidate generation.

Input:

- D , a database of transactions;
- min_sup , the minimum support count threshold.

Output: L , frequent itemsets in D .

Method:

```

(1)  $L_1 = \text{find\_frequent\_1-itemsets}(D)$ ;
(2) for  $(k = 2; L_{k-1} \neq \emptyset; k++)$  {
(3)    $C_k = \text{apriori\_gen}(L_{k-1})$ ;
(4)   for each transaction  $t \in D$  { // scan  $D$  for counts
(5)      $C_t = \text{subset}(C_k, t)$ ; // get the subsets of  $t$  that are candidates
(6)     for each candidate  $c \in C_t$ 
(7)        $c.\text{count}++$ ;
(8)   }
(9)    $L_k = \{c \in C_k | c.\text{count} \geq min\_sup\}$ 
(10) }
(11) return  $L = \cup_k L_k$ ;

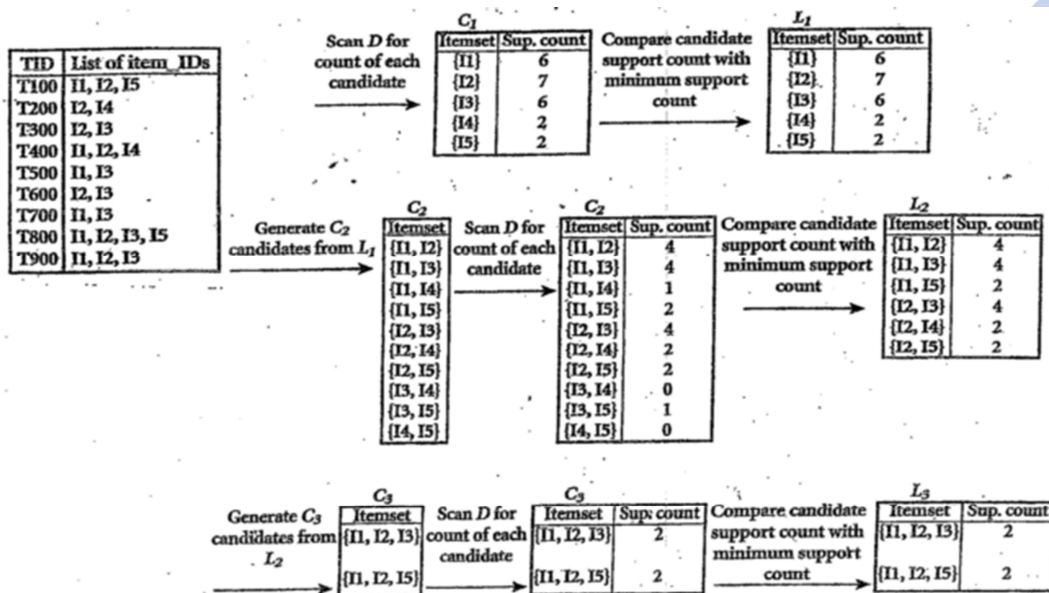
procedure apriori_gen( $L_{k-1}$ : frequent  $(k-1)$ -itemsets)
(1) for each itemset  $l_1 \in L_{k-1}$ 
(2)   for each itemset  $l_2 \in L_{k-1}$ 
(3)     if  $(l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2]) \wedge \dots \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$  then {
(4)        $c = l_1 \bowtie l_2$ ; // join step: generate candidates
(5)       if  $\text{has\_infrequent\_subset}(c, L_{k-1})$  then
(6)         delete  $c$ ; // prune step: remove unfruitful candidate
(7)       else add  $c$  to  $C_k$ ;
(8)     }
(9) return  $C_k$ ;

procedure has_infrequent_subset( $c$ : candidate  $k$ -itemset;
                                 $L_{k-1}$ : frequent  $(k-1)$ -itemsets; // use prior knowledge
(1) for each  $(k-1)$ -subset  $s$  of  $c$ 
(2)   if  $s \notin L_{k-1}$  then
(3)     return TRUE;
(4) return FALSE;

```

2110773-3 2/67

17



2110773-3 2/67

18