

Association Mining

วัตถุประสงค์ของการทำ Association Mining คือการค้นหาคำสัมพันธ์ที่น่าสนใจ เช่น
ลูกค้าที่ซื้อผ้าอ้อมส่วนใหญ่ จะซื้อ เปียม ด้วย

เขียนกฎความสัมพันธ์ได้เป็น

$\{\text{ผ้าอ้อม}\} \longrightarrow \{\text{เปียม}\}$

$\{x\} \longrightarrow \{y\}$

$\{a,b,c\} \longrightarrow \{d\}$

$\{v,w\} \longrightarrow \{x,y,z\}$

ในข้อมูล transaction จริงๆ จะมีรายการสินค้าที่ถูกขายมากมาย ซึ่งการจับคู่ หรือจับกลุ่ม

รายการสินค้าเหล่านี้ทั้งหมด เพื่อหากฎความสัมพันธ์
เป็นความคิดที่ไม่ค่อยดีนัก เพราะ $n = \text{จำนวนประเภทข้อมูลทั้งหมด}$
ขนาด $(2^n - n - 1)$; $n = \text{จำนวนรายการสินค้าทั้งหมด}$

กฎทั้งหมดที่ได้มาส่วนใหญ่ จะมีค่า Confidence $< \text{minimum Confidence}$

★ เพื่อหลีกเลี่ยงการประมวลผลข้อมูลจำนวนมาก ★

เราจำเป็นต้องคัดกรองการสินค้าที่ ปรากฏบ่อย ก่อน แล้วจึงสร้างกฎความสัมพันธ์

Q: ทำไมต้องคัดกรอง = สินค้าที่ปรากฏบ่อย ?

A: เพราะการหาความสัมพันธ์ของสินค้าที่มีการขายน้อยครั้ง จะให้กฎความสัมพันธ์ ที่เป็น noise
(กฎความสัมพันธ์ใช้งานจริงไม่ได้)

ขั้นตอน Association Mining

1. หาเซตของรายการสินค้าที่ปรากฏบ่อย \longrightarrow 2. หากฎความสัมพันธ์บนเซตเหล่านั้น

① - How to หาเซตของรายการสินค้าที่ปรากฏบ่อย

Q: เราจะสามารถรู้ได้ว่า เซตของรายการสินค้าปรากฏบ่อย หรือไม่น้อย ?

A: เซตของรายการสินค้า ที่มีค่า Support $\geq \text{minimum Support}$

จะเป็น เซตของรายการสินค้าที่ปรากฏบ่อย

ค่า Support สามารถหาได้จาก

$$\text{support}(x \rightarrow y) = \frac{\text{จำนวน transaction ที่มี } x \text{ และ } y \text{ ปรากฏร่วมกัน}}{\text{จำนวน transaction ทั้งหมด}}$$

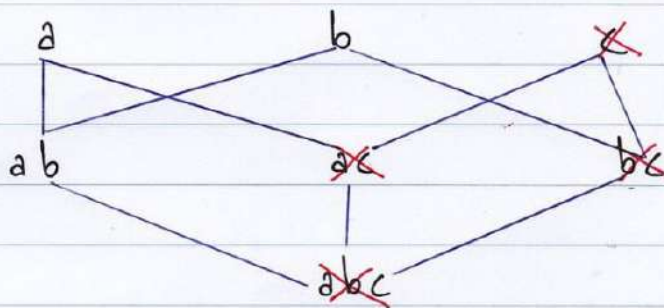
เขียนเป็นสูตร ดังนี้
$$\text{support}(x \rightarrow y) = \frac{\text{tran_count}(x \cup y)}{\text{tran_count}(D)}$$

* minimum support เป็นค่าที่เราต้องเลือกเอง

$$\text{โดยที่ } 0 \leq \text{minimum support} \leq 1$$

การหาเซตทั้งหมดของรายการสินค้าที่ปรากฏบ่อยไม่ใช่ความคิดที่ดีนัก ที่จริงจากจะคำนวณประมวลผลข้อมูลจำนวนมหาศาล ($2^n - 1$)

วิธีการแก้ก็คือ เราไม่ต้องสร้างเซตของรายการสินค้าที่มีสมาชิก เป็นเซตที่ไม่ปรากฏบ่อย



จากรูปนี้ เซตของสินค้า c เป็นเซตที่มีค่า support น้อยกว่า minimum support เมื่อ $\text{support}(c) < \text{minimum support}$ แสดงว่า เซตของสินค้า c เป็นเซตที่ ตัดทิ้ง ดังนั้น เซตของ $a c$, $b c$ และ $a b c$ ก็จะเป็นเซตที่ปรากฏไม่บ่อยด้วย ปรากฏไม่บ่อย จึงไม่จำเป็นต้องสร้างเซตเหล่านี้ขึ้นมา

วิธีการนี้มีชื่อว่า APRIORI Algorithm

Example of APRIORI Algorithm

Tran-ID	List of Item
T1	I1, I2, I5
T2	I2, I4
T3	I2, I3
T4	I1, I3, I4
T5	I1, I3
T6	I2, I3
T7	I1, I3
T8	I1, I2, I3, I5
T9	I1, I3, I3

ขั้นตอนที่ 1

นับจำนวนรายการสินค้าทั้งหมด

ในที่นี้จะได้ว่า มี รายการสินค้าทั้งหมด 5 รายการ

ได้แก่ I1, I2, I3, I4, I5

ขั้นตอนที่ 2

สร้าง Candidate itemset (C_1)

ทำการนับว่ารายการสินค้าแต่ละตัวในขั้นตอนที่ 1 ปรากฏกี่ครั้ง
จะได้ว่า { ใน transaction

I1 ปรากฏ 6 ครั้ง

I2 ปรากฏ 7 ครั้ง

I3 ปรากฏ 6 ครั้ง

I4 ปรากฏ 2 ครั้ง

I5 ปรากฏ 2 ครั้ง

ต่อมาจะตรวจสอบค่านี้ไปสร้าง Candidate itemset

Candidate itemset ประกอบด้วย

ซึ่ง itemset นี้คือ ถ้า support ของ itemset นั้น
ไม่ติด (จำนวน transaction ทั้งหมด = 9)

Itemset	Support
$\{I1\}$	6/9
$\{I2\}$	7/9
$\{I3\}$	6/9
$\{I4\}$	2/9
$\{I5\}$	2/9

$C_1 =$

$= L_1$

ขั้นตอนที่ 3

สร้าง frequent itemset (L_2)

โดยการพิจารณาแต่ละ itemset ใน C_1 ถ้ามีค่า support (itemset) \geq minimum support ใน

ถ้าค่า support ของแต่ละ itemset ใน $C_1 \geq$ minimum support

itemset นั้นก็จะผ่านการคัดเลือก ไม่ใช้สร้าง frequent itemset (ในที่นี้ minimum support

ขั้นต่ำ ทุก itemset ใน C_1 จึงผ่านการคัดเลือก ไม่ใช้สร้าง L_1 = 2/9)

ขั้นตอนที่ 4

สร้าง itemset C_2

เราสามารถสร้าง itemset C_2 ได้จาก frequent itemset (L_1)

เพราะว่า itemset $L_1 = \{I1\}, \{I2\}, \{I3\}, \{I4\}, \{I5\}$

ดังนั้นจะสามารถสร้าง itemset C_2 ได้ทั้งหมด 10 แบบ ได้แก่

$\{I1, I2\}, \{I1, I3\}, \{I1, I4\}, \{I1, I5\}, \{I2, I3\}, \{I2, I4\}, \{I2, I5\}, \{I3, I4\}, \{I3, I5\}, \{I4, I5\}$

แต่ในการใช้งานจริง เราไม่ต้องสร้าง Itemset C_2 ขึ้นมาครบทุกแบบตามที่เขียนด้านบน

เราทำการสร้างเฉพาะ itemset ที่มีการปรากฏใน transaction ก็เพียงพอ (เพื่อลดขั้นตอนในการ
จากการพิจารณาพบว่า - การประมวลผลในลำดับถัดไป)

$\{I1, I2\}$	ปรากฏใน	T1, T4, T8, T9
$\{I1, I3\}$	"	T5, T7, T8, T9
$\{I1, I4\}$	"	T4
$\{I1, I5\}$	"	T1, T8
$\{I2, I3\}$	"	T3, T8, T9
$\{I2, I4\}$	"	T3, T4
$\{I2, I5\}$	"	T1, T8
$\{I3, I4\}$	"	ไม่ปรากฏใน transaction ใด
$\{I3, I5\}$	"	T8
$\{I4, I5\}$	"	ไม่ปรากฏใน transaction ใด

จากข้อมูลด้านบนจะเห็นว่า $\{I3, I4\}, \{I4, I5\}$ ไม่มีการปรากฏใน transaction ใดๆ
จะได้ว่า

Itemset

$\{I1, I2\}$

$\{I1, I3\}$

$\{I1, I4\}$

$\{I1, I5\}$

$\{I2, I3\}$

$\{I2, I4\}$

$\{I2, I5\}$

$\{I3, I5\}$

*** จะเห็นว่า การสร้าง Itemset C_2

ในขั้นตอนนี้ เป็นเพียงการพิจารณาเท่านั้น

ว่าต้องสร้าง itemset ใดขึ้นมาบ้าง

(เป็นประโยชน์ต่อการลดขั้นตอนในลำดับถัดไป) ***

ขั้นตอนที่ 5 (หลังจากขั้นตอนที่ 2)

สร้าง (candidate itemset (C_2))

ทำการนับว่า itemset แต่ละตัวในขั้นตอนที่ 4 ปรากฏกี่ครั้งใน transaction

q=9 ตัวอย่าง

{I1, I2}	ปรากฏ	4	ครั้ง
{I1, I3}	"	4	ครั้ง
{I1, I4}	"	1	ครั้ง
{I1, I5}	"	2	ครั้ง
{I2, I3}	"	3	ครั้ง
{I3, I4}	"	2	ครั้ง
{I3, I5}	"	2	ครั้ง
{I3, I5}	"	1	ครั้ง

พิจารณาข้อมูล เหล่านี้ ไม่สร้าง Candidate itemset

นี่คือ

	Itemset	Support
	{I1, I2}	4/9
	{I1, I3}	4/9
	{I1, I4}	1/9
$C_2 =$	{I1, I5}	2/9
	{I3, I3}	3/9
	{I3, I4}	2/9
	{I3, I5}	2/9
	{I3, I5}	1/9

ขั้นตอนที่ 6 (เหมือนขั้นตอนที่ 3)สร้าง frequent itemset (L_2)

โดยการพิจารณา แต่ละ itemset ใน C_2 ว่ามีค่า $\text{support}(\text{itemset}) \geq \text{minimum support}$ หรือไม่
 จากการพิจารณาจะได้ว่า (ในที่นี้ $\text{minimum support} = 2/9$)

$L_2 =$	Itemset	Support
	{I1, I2}	4/9
	{I1, I3}	4/9
	{I1, I5}	2/9
	{I3, I3}	4/9
	{I3, I4}	2/9
	{I3, I5}	2/9

ขั้นตอนที่ 7 (เหมือนขั้นตอนที่ 4)สร้าง itemset C_3 เราสามารถสร้าง itemset C_3 ได้จาก frequent itemset (L_2)เพราะว่า itemset $L_2 = \{I1, I2\}, \{I1, I3\}, \{I1, I5\}, \{I3, I3\}, \{I3, I4\}, \{I3, I5\}$ ดังนั้นเราสามารถสร้าง itemset C_3 ได้ทั้งหมด 4 แบบ ได้แก่

{I1, I2, I3}, {I1, I3, I5}, {I3, I3, I4}, {I3, I3, I5}

แต่ในการใช้งานจริง เราไม่ต้องการสร้าง itemset C_3 ขึ้นมาครบทุกแบบตามที่เขียนด้านบน

เราทำการสร้างเฉพาะ itemset ที่มีการปรากฏใน transaction ก็พอ

จากการพิจารณาพบว่า

{I1, I3, I3}	ปรากฏใน	T8, T9
{I1, I3, I5}	ปรากฏใน	T1, T9
{I3, I3, I4}	ไม่ปรากฏใน transaction ใด	
{I3, I3, I5}	ไม่ปรากฏใน transaction ใด	

จากข้อมูลด้านบนจะเห็นว่า {I3, I3, I4}, {I3, I3, I5} ไม่มีการปรากฏใน transaction ใดๆ
 จึงได้ว่า

Itemset
{I1, I2, I3}
{I1, I3, I5}

ขั้นตอนที่ 8 (เหมือนขั้นตอนที่ 2, 5)

สร้าง Candidate itemset (C_3)

ทำการเพิ่ม itemset แต่ละตัวในขั้นตอนที่ 7 ปรากฏก็จริงใน transaction
จะได้ว่า

$\{I1, I2, I3\}$ ปรากฏ 2 ครั้ง

$\{I1, I2, I5\}$ ปรากฏ 2 ครั้ง

ผลมาจากการนับค่านี้จะได้สร้าง Candidate itemset
ขึ้นดังนี้

Itemset	Support
$C_3 = \{I1, I2, I3\}$	2/9
$\{I1, I2, I5\}$	2/9

ขั้นตอนที่ 9 (เหมือนขั้นตอนที่ 3, 6)

สร้าง frequent itemset (L_3)

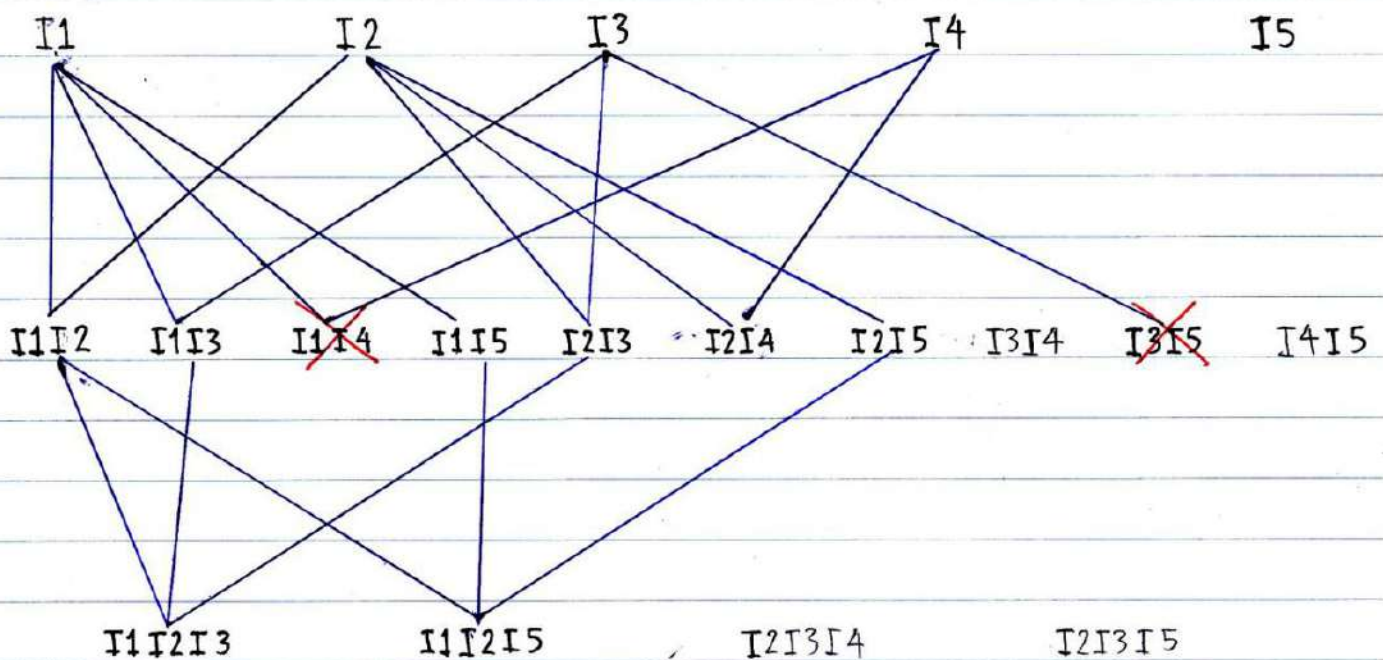
โดยการพิจารณาแต่ละ itemset ใน C_3 ว่ามีค่า $\text{support}(\text{itemset}) \geq \text{minimum support}$ หรือไม่
จากการพิจารณาจะได้ว่า (ในที่นี้ $\text{minimum support} = 2/9$)

Itemset	Support
$L_3 = \{I1, I2, I3\}$	2/9
$\{I1, I2, I5\}$	2/9

ถ้าเราพิจารณา itemset C_4 เราพบว่า itemset $C_4 = \{I1, I2, I3, I5\}$
ซึ่งไม่มีปรากฏใน transaction ใดๆ

ดังนั้น

“
END
”



รูปแสดงการสร้าง frequent itemset
ด้วย APRIORI Algorithm

หลังจากที่เราหาเซตของรายการสินค้าที่ปรากฏบ่อย (frequent itemset) เสร็จแล้ว
ต่อมาเราก็ทำการหา กฎ ความสัมพันธ์ ได้โดย

- ②
- How to หา กฎความสัมพันธ์
- Q: เราจะทราบได้อย่างไรว่า เซตของรายการสินค้ามีความสัมพันธ์กัน ?
- A: เซตของรายการสินค้าที่มีค่า $\text{Confidence} \geq \text{minimum Confidence}$
จะเป็นเซตของรายการสินค้าที่มีความสัมพันธ์กัน

ค่า Confidence สามารถหาได้จาก

$$\text{confidence}(x \rightarrow y) = \frac{\text{จำนวน transaction ที่มี } x \text{ และ } y \text{ ปรากฏร่วมกัน}}{\text{จำนวน transaction ที่มี } x \text{ ปรากฏ}}$$

เขียนเป็นสูตรได้ดังนี้ $\text{confidence}(x \rightarrow y) = \frac{\text{tran_count}(x \cup y)}{\text{tran_count}(x)}$

* minimum confidence เป็นค่าที่เราตั้งไว้

$$\text{โดยที่ } 0 \leq \text{minimum confidence} \leq 1$$

เราเรียก $\text{confidence}(x \rightarrow y)$ ที่ $\geq \text{minimum confidence}$

ว่า "Strong Association Rule"

Example of Confidence Calculation

เราใช้ frequent itemset ที่ได้จากขั้นตอนที่ 1 (ข้างบนที่แล้ว)
มาหา confidence

เราไม่สนใจ frequent itemset ที่มีสมาชิกเพียง 1 ตัว ในตารางความสัมพันธ์
จะได้ว่า

$$\text{confidence}(I1 \rightarrow I2) = \frac{\text{tran_count}(I1 \cup I2)}{\text{tran_count}(I1)} = \frac{4}{6}$$

$$\text{confidence}(I2 \rightarrow I1) = \frac{\text{tran_count}(I1 \cup I2)}{\text{tran_count}(I2)} = \frac{4}{7}$$

* จากตาราง เราสังเกตได้ว่า $\text{confidence}(x \rightarrow y)$ ไม่จำเป็นต้อง
เท่ากับ $\text{confidence}(y \rightarrow x)$

$$\text{confidence}(I1, I2 \rightarrow I3) = \frac{\text{tran_count}(I1 \cup I2 \cup I3)}{\text{tran_count}(I1 \cup I2)} = \frac{2}{4}$$

$$\text{confidence}(I1 \rightarrow I2, I3) = \frac{\text{tran_count}(I1 \cup I2 \cup I3)}{\text{tran_count}(I1)} = \frac{2}{6}$$

$$\text{confidence}(I3 \rightarrow I1, I2) = \frac{\text{tran_count}(I1 \cup I2 \cup I3)}{\text{tran_count}(I3)} = \frac{2}{6}$$

$$\text{confidence}(I1, I3 \rightarrow I2) = \frac{\text{tran_count}(I1 \cup I2 \cup I3)}{\text{tran_count}(I1 \cup I3)} = \frac{2}{4}$$

$$\text{confidence}(I2 \rightarrow I1, I3) = \frac{\text{tran_count}(I1 \cup I2 \cup I3)}{\text{tran_count}(I2)} = \frac{2}{7}$$

$$\text{confidence}(I2, I3 \rightarrow I1) = \frac{\text{tran_count}(I1 \cup I2 \cup I3)}{\text{tran_count}(I2 \cup I3)} = \frac{2}{3}$$

* ขั้นตอนสุดท้าย หลังจาก ได้ กฎความสัมพันธ์ (Strong Association Rule) แล้ว ก็คิด
การทดสอบว่า Strong Association Rule สามารถนำไปใช้
งานจริงได้ไหม?

ยกตัวอย่างเช่น

$$\text{confidence}(A \rightarrow B) = 0.7 \quad \text{นั่นหมายความว่า มีโอกาส 70\%}$$

ที่ลูกค้าซื้อสินค้า A แล้วจะซื้อสินค้า B ด้วย

แต่ถ้า เราทราบว่า ความน่าจะเป็นที่ลูกค้าซื้อสินค้า B เป็น 90% ของจำนวน transaction
แล้ว Strong Association Rule ที่เราได้มา น่าจะดีกว่าไหม? { หักลบ }

③

- วิธีการตรวจสอบว่ากฎความสัมพันธ์ที่ได้มา สามารถนำไปใช้ได้จริงไหม มีดังนี้

① หาค่าสนับสนุน ซึ่งมีสูตร $\text{corr}_{A,B} = \frac{P(A \cup B)}{P(A)P(B)}$

② พิจารณาค่า $\text{corr}_{A,B}$ ที่ได้

- ถ้า $\text{corr}_{A,B} \geq 1$ แล้ว Strong Association Rule สามารถใช้งานได้จริง

- แต่ถ้า $\text{corr}_{A,B} < 1$ แล้ว Strong Association Rule ไม่สามารถใช้งานได้จริง
(การซื้อสินค้าชนิดหนึ่ง จะลด โอกาสที่จะซื้อ -
สินค้าอีกชนิดหนึ่ง)