

(Unsupervised Learning)

Association Rule Mining

Examples

→ Market Basket Analysis
→ Recommendation Systems

Sale

→ Web log analysis

Frequent pattern

→ pattern occurs frequently in a dataset.

Generate association rules from frequent patterns.

Dataset

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

$\frac{n}{N}$

* Support Count (S)

→ Frequency of occurrence of an itemset

→ e.g. in above example {Bread, Milk, Diaper} occurs two times

* Support

→ Fraction of transactions that contain an itemset

{Milk, Bread, Diaper}

Formula:

$$\frac{n}{N} = \frac{2}{5}$$

support will tell frequent item

Association Rule

- 1) $X \rightarrow Y$ (support)
- 2) $Y \rightarrow X$ (support)

X and Y are itemsets

For Association $\{Milk, Diaper\} \rightarrow Beer$

we need

Rule Evaluation Metrics

Support (s) : $\frac{n}{N}$

Confidence : how often items in Y appear in transactions that contain X

$$A \rightarrow B \Rightarrow P(B/A) = \frac{P(A \cup B)}{P(A)} \quad \text{or} \quad \frac{\sigma(A, B)}{\sigma(A)}$$

if we want to find support and confidence for this case

$\{Milk, Diaper\} \Rightarrow Beer$

$s = \frac{2}{5}$ As it occurs 2 times

$c = \frac{\frac{2}{5}}{\frac{3}{5}} = \frac{2}{3}$

$c = \frac{\sigma(Milk, Diaper, Beer)}{\sigma(Milk, Diaper)}$

find all rules

support \geq minsup threshold

confidence \geq minconf threshold

Brute-force approach

S becomes

$\{A, B, C\}$

$\{A\} \rightarrow \{B, C\}, \{B\} \rightarrow \{A, C\}, \{C\} \rightarrow \{A, B\}$

we is

all set

and

we

which contains only

$\{A, B\} \rightarrow \{C\}, \{B, C\} \rightarrow \{A\}, \{A, C\} \rightarrow \{B\}$

$$2^3 = 8$$

Apriori Algorithm

All subsets must be frequent

Anti-monotone property

$$(\forall X \subseteq Y) \quad s(X) \geq s(Y)$$

Pruning

The superset of an infrequent itemset
are infrequent

two operations join \rightarrow prune

C_k : Candidate itemset of size k

L_k : Frequent itemset of size k



$L_1 = \{ \text{frequent items} \}$

for $(k=1; L_k \neq \emptyset; k++)$
do begin

join steps: join L_k with itself to produce C_{k+1}

PRUNE steps: discard $(k+1)$ -itemsets from C_{k+1} that contain non-frequent k -itemsets as subsets

$C_{k+1} =$ candidates generated from L_k

for each transaction t in database do
increment the count of all
candidates in C_{k+1} that are contained
in t
 L_{k+1} is candidates in C_{k+1} with
min-support

end

return $\bigcup_k L_k$

union of all those

two items
must be
common for joining

Step 1:

self-joining L_k

Joining $(C_{k-1} \times L_{k-1} \times L_{k-1})$

The join, $L_{k-1} \times L_{k-1}$, is performed only if their first $(k-2)$ items are in common.

Example of Candidate generation

$L_2 = \{abc, abd, acd, ace, bcd\}$

Self-joining: $L_2 \times L_2$

~~abc~~ ~~abd~~ ~~acd~~ ~~ace~~ ~~bcd~~
two items
abc, abd
acd, ace

Pruning

abcd

↓

abc, abd, bcd,
all subsets are
present in L_3

acde X

↓

acd, ade X not
include

$C_4 = \{abcd\}$

A priori Example

Transaction ID	Items
T1	{ I1, I3, I4, I6 }
T2	{ I2, I3, I5 }
T3	{ I1, I2, I3, I6 }
T4	{ I1, I5, I6 }

Iteration 1

Iteration 1

L_1

Itemset	Count
$\{1\}$	3
$\{1, 2\}$	2
$\{1, 3\}$	3
$\{1, 4\}$	1
$\{2, 5\}$	3
$\{2, 6\}$	2

L_2

Itemset	Count
$\{1, 2\}$	3
$\{1, 3\}$	2
$\{1, 4\}$	3
$\{2, 5\}$	3
$\{2, 6\}$	2

0.5

0.5

100

4

Min support threshold 0.5

$$L_1 \otimes L_1$$

$$\begin{array}{r} 50 \\ \times 4 \\ \hline 200 \end{array}$$

1-2 element
comms

Iteration 2

L_1	X	L_1	$K-2$	$2-2=0$
Itemsel		Item	$4-2=2$	
$\{I_1, I_2\}$		$\{I_1\}$		
$\{I_2\}$		$\{I_2\}$		
$\{I_5\}$		$\{I_3\}$		$\{I_1, I_2\}$
$\{I_5\}$		$\{I_5\}$		
$\{I_6\}$		$\{I_6\}$		$\{I_1, I_2\}$

b

C_2	$\{I_1, I_2\}$	1	X
	$\{I_1, I_3\}$	2	
	$\{I_1, I_5\}$	2	
	$\{I_1, I_6\}$	2	
	$\{I_2, I_3\}$	2	
	$\{I_2, I_5\}$	2	
	$\{I_2, I_6\}$	2	
	$\{I_3, I_5\}$	2	X
	$\{I_3, I_6\}$	1	X
	$\{I_5, I_6\}$	1	X

L_2

$\{I_1, I_3\}$	2
$\{I_1, I_5\}$	2
$\{I_1, I_6\}$	2
$\{I_2, I_3\}$	2
$\{I_2, I_5\}$	2
$\{I_3, I_5\}$	2

$L_2 \times L_2$

1 element
should be
common

$\{I_1, I_3\}$	$\{I_1, I_3\}$
$\{I_1, I_5\}$	$\{I_1, I_5\}$
$\{I_1, I_6\}$	$\{I_1, I_6\}$
$\rightarrow \{I_2, I_3\}$	$\{I_2, I_3\}$
$\{I_2, I_5\}$	$\{I_2, I_5\}$
$\{I_3, I_5\}$	$\{I_3, I_5\}$

Q3

$\{I_1, I_3, I_5\}$	
$\{I_1, I_3\}$	
$\{I_1, I_3, I_5\}$	\neq
$\{I_1, I_3, I_6\}$	\neq
$\{I_2, I_3, I_5\}$	
$\{I_1, I_5, I_6\}$	\neq
$\{I_2, I_3, I_6\}$	

$\checkmark \rightarrow C_3$	$\{I_1, I_3, I_5\}$	$\{I_1, I_3\}, \{I_1, I_5\}, \{I_3, I_5\}$	\checkmark
$\rightarrow X$	$\{I_1, I_3, I_6\}$	$\{I_1, I_3\}, \{I_1, I_6\}, \{I_3, I_6\}$	\times
\rightarrow	$\{I_1, I_5, I_6\}$	$\{I_1, I_5\}, \{I_1, I_6\}, \{I_5, I_6\}$	\times
$\checkmark \rightarrow$	$\{I_2, I_3, I_6\}$	$\{I_2, I_3\}, \{I_2, I_6\}, \{I_3, I_6\}$	\checkmark

Q3

~~$\{I_1, I_3, I_5\}$~~ , ~~$\{I_2, I_3, I_5\}$~~

C_3	$\{I_1, I_2, I_3\}$	1
	$\{I_2, I_3, I_5\}$	2 \checkmark

L_3

$$\{I_2, I_3, I_5\} \quad | \quad 2$$

two sets need to join

C_1 is empty

Frequent Item set $\{I_2, I_3, I_5\}$

L_1		L_2		L_3	
$\{I_1\}$	3	$\{I_1, I_3\}$	2	$\{I_2, I_3, I_5\}$	1
$\{I_2\}$	2	$\{I_1, I_5\}$	2		
$\{I_3\}$	3	$\{I_1, I_6\}$	2	$\{I_2, I_3, I_5\}$	
$\{I_5\}$	3	$\{I_2, I_3\}$	2	$\{I_2, I_3\} \rightarrow \{I_5\}$	
$\{I_6\}$	2	$\{I_2, I_5\}$	2		
		$\{I_3, I_5\}$	2		

For L_3

$$\begin{aligned} \{I_2, I_3\} &\rightarrow \{I_5\} \\ \{I_3, I_5\} &\rightarrow \{I_2\} \\ \{I_2, I_5\} &\rightarrow \{I_3\} \\ \{I_2\} &\rightarrow \{I_3, I_5\} \\ \{I_3\} &\rightarrow \{I_2, I_5\} \\ \{I_5\} &\rightarrow \{I_2, I_3\} \end{aligned}$$

$$\frac{P(I_2 \cup I_3 \cup I_5)}{P(I_2 \cup I_3)}$$

2/2
2/2
2/2
2/2
2/3
2/3

60

Example 2.

Transaction ID	Items
T100	{M, O, N, K, E, Y}
T200	{D, O, N, K, E, Y}
T300	{M, A, K, E}
T400	{M, U, C, K, Y}
T500	{C, O, K, E}

Transaction ID	Items
T100	{E, K, M, N, O, Y}
T200	{D, E, K, N, O, Y}
T300	{A, E, K, M}
T400	{C, K, M, U, Y}
T500	{C, E, I, K, O, Q}

9

Itemset	Count
{A}	1
{C}	2
{D}	1
{E}	4
{Z}	1
{K}	5
{M}	3
{N}	2
{O}	3
{U}	2
{Y}	3

Min Support

60%

$$= \frac{60}{100} \times 100$$

3 3

L₁

Itemset	Count
{G}	4
{K}	5
{M}	3
{O}	3
{Y}	3

L₁

Itemset	Item
{E}	{E}
{K}	{K}
{M}	{M}
{O}	{O}
{Y}	{Y}

L₁

{E, K}

C ₂	Itemset	
	{E, K}	4
	{E, M}	2
	{E, O}	3
	{E, Y}	2
	{K, M}	3
	{K, O}	3
	{K, Y}	3
	{M, O}	1
	{M, Y}	2
	{O, Y}	2

L₂ ~~6~~

Itemset	Count
{E, K}	4
{E, O}	3
{K, M}	3
{K, O}	3
{K, Y}	3
{E, M}	1

L₂ x L₂

{E, K}	{E, K}
{E, O}	{E, O}
{K, M}	{K, M}
{K, O}	{K, O}
{K, Y}	{K, Y}

~~Itemset~~
 ~~$\{E, k, 0\}$~~
 ~~$\{E, k, M\}$~~
 ~~$\{E, k, Y\}$~~

C₂

Itemset

$\{E, k, 0\}$

$\{k, M, 0\}$

$\{k, M, Y\}$

$\{k, 0, X\}$

Itemset

$\{E, k, 0\}$

$\{E, k\}, \{k, 0\}, \{E, 0\}$ ✓

$\{k, M, 0\}$

$\{k, M\}, \{k, 0\}, \{M, 0\}$ X

$\{k, M, Y\}$

$\{k, M\}, \{M, Y\}, \{k, Y\}$ X

$\{k, 0, Y\}$

$\{k, 0\}, \{0, Y\}, \{k, Y\}$ X

C₃

$\{E, k, 0\}$

3

L₃

$\{E, k, 0\}$

3

C₁ is empty

C₄ is empty

L₁ L₂ L₃

Itemset	Count
{E}	4
{K}	5
{M}	3
{O}	3
{Y}	3

L₂

Itemset	Count
{E, K}	4
{E, O}	3
{K, M}	3
{K, O}	3
{K, Y}	3

L₃

Itemset	Count
{E, K, O}	3

{E, K} → {O}	$\frac{3}{4}$	= 75%
{K, O} → {E}	$\frac{3}{3}$	= 100%
{E, O} → {K}	$\frac{3}{3}$	= 100%
{O} → {E, K}	$\frac{3}{3}$	= 100%
{E} → {K, O}	$\frac{3}{4}$	= 75%
{K} → {E, O}	$\frac{3}{5}$	= 60%

FP Growth Algorithm

Frequent Pattern

→ tree

TID	List of item-IDs
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3, I3 , I5
T800	I1, I2, I3, I5
T900	I1, I2, I3

We have to scan database once

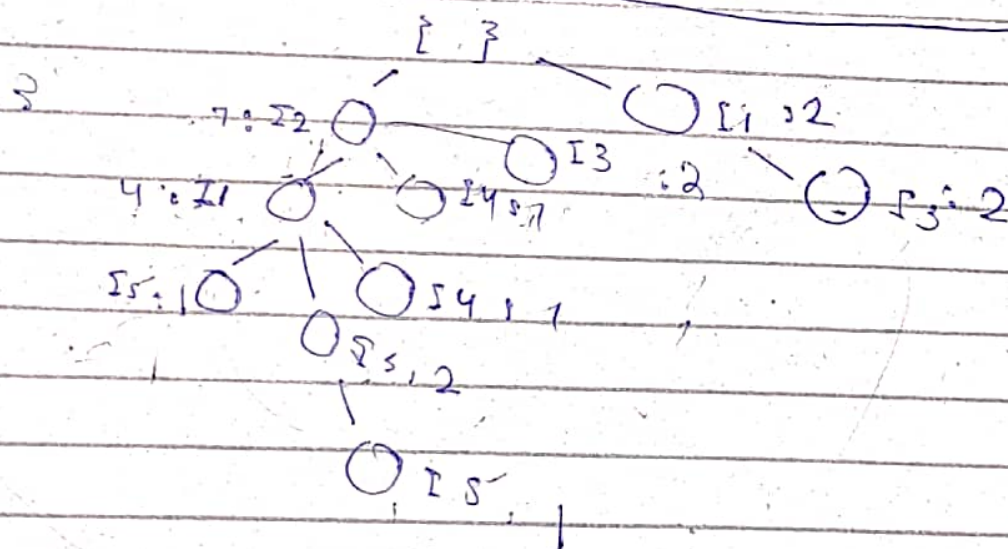
Itemset	Count	order	Itemset	Count
{I1}	6	→	{I2}	7
{I2}	7		{I1}	6
{I3}	6		{I3}	6
{I4}	2		{I4}	2
{I5}	2		{I5}	2

Construct FP tree

arrange the table having more cont

Arranged Table

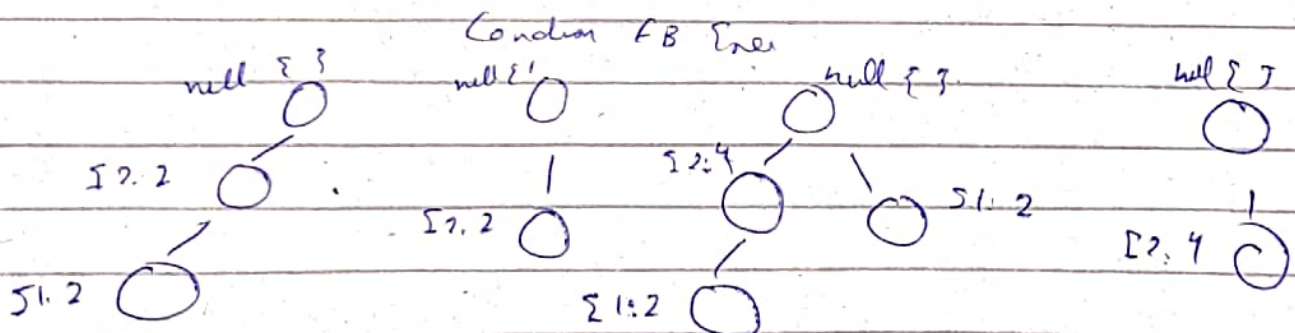
Itemset	Count
T100	I2, I1, I5
T200	I2, I4
T300	I2, I3
T400	I2, I1, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I2, I1, I3, I5
T900	I2, I1, I3



ItemID	Count
{I5}	2
{I4}	2
{I3}	6
{I1}	6
{I2}	7

ItemID	Conditional Pattern Base
{I5}	{I2, I1: 1}, {I2, I1, I3: 1}
{I4}	{I2, I1: 1}, {I2, I1}
{I3}	{I2, I1, I2}, {I2, I1, I2}, {I1: 2}
{I1}	{I2: 4}

ItemID	Conditional FP Tree
{I5}	<I2: 2, I1: 2>
{I4}	<I2: 2>
{I3}	<I2: 4, I1: 2> <I1: 2>
{I1}	<I2: 4>



Item ID	Frequent Pattern
$\{15\}$	$\{12, 15, 2\}, \{11, 15, 2\}, \{12, 11, 15, 2\}$
$\{24\}$	$\{12, 24, 2\}$
$\{13\}$	$\{12, 13, 4\}, \{11, 13, 4\}, \{12, 11, 13, 2\}$
$\{11\}$	$\{12, 11, 4\}$
$\{52\}$	

Frequent Pattern

$\{11\}, \{12\}, \{13\}, \{14\}, \{15\}, \{12, 15\},$
 $\{11, 15\}, \{12, 14\}, \{12, 13\}, \{11, 13\},$
 $\{12, 11, 15\}, \{12, 11, 13\}$

Example 2

TID	items bought
T100	{M, O, N, K, E, Y}
T200	{D, O, N, K, E, Y}
T300	{M, A, K, E}
T400	{M, U, C, K, Y}
T500	{C, O, O, K, E}

min sup = 60%

$$\begin{aligned}
 &= 60\% \\
 &= 60/100 \times 5 \\
 &= 3
 \end{aligned}$$

Itemset	Count
{M}	3
{O}	3
{N}	2
{K}	5
{E}	4
{Y}	3
{D}	1
{A}	1
{U}	1
{C}	2
{I}	1

Itemset	Count
{M}	3
{O}	3
{K}	5
{E}	4
{Y}	3

Itemset	Count
{K}	5
{E}	4
{M}	3
{O}	3
{Y}	3

Transaction ID

Items

T100

$\{k, e, m, o, y\}$

T200

$\{k, e, o, y\}$

T300

$\{k, e, m\}$

T400

$\{k, m, y\}$

T500

$\{k, e, o\}$

$\{ \}$ null

k

M: 1

e: 4

y: 1

o: 2

M: 3

y: 1

o: 1

y: 1

Item ID

Conditional Pattern Base

$\{y\}$

$\{ \{k, e, m, o, y\} : 1\}, \{ \{k, e, o, y\} : 1\}, \{ \{k, m, y\} : 1\}$

$\{o\}$

$\{ \{k, e, m\} : 1\}, \{ \{k, e\} : 2\}$

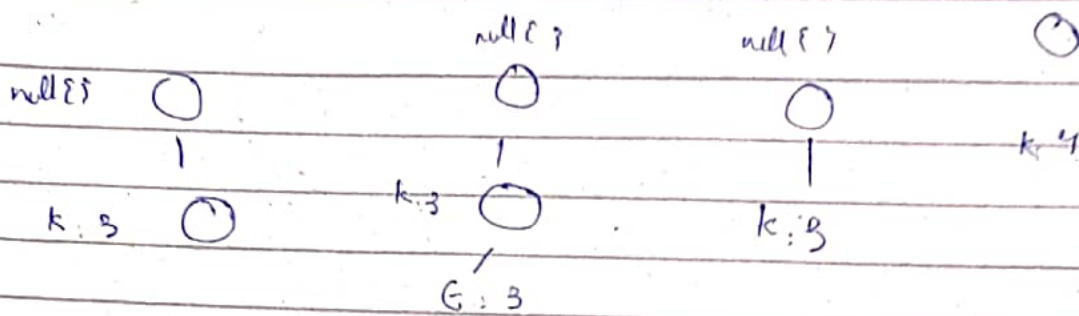
$\{m\}$

$\{ \{k, e\} : 3\}, \{ \{k, m\} : 1\}$

$\{e\}$

$\{ \{k\} : 4\}$

$\{k\}$



Item ID	Pattern Base	Conditional FP Tree
$\{Y\}$	$\{K, E, M, O, Y: 1\}, \{K, G, O: 1\}, \{K, M: 1\}$	$\{K: 3\}$
$\{O\}$	$\{K, G, M: 1\}, \{K, E: 2\}$	$\{K, E: 3\}$
$\{M\}$	$\{K, E: 2\}, \{K: 1\}$	$\{K: 3\}$
$\{E\}$	$\{K: 4\}$	$\{K: 4\}$
$\{K\}$		

Frequent Patterns

- $\{KY: 3\}$
- $\{KO: 3, EO: 3, KEO: 3\}$
- $\{KM: 3\}$
- $\{KE: 4\}$

Frequent

- $\{Y\}$ $\{O\}$ $\{M\}$ $\{E\}$ $\{K\}$
- $\{KY\}$ $\{KO\}$ $\{EO\}$ $\{KEO\}$ $\{KM\}$
- $\{KE\}$