

1) a) star schema diagram :

time dimension table

time-key
day
month
quarter
Year

Patient dimension table

Patient-key
Patient-name
Phone-number
Address
Description

Fact table

time-key
doctor-key
Patient-key
charge
count

doctor dimension table

doctor-key
doctor-name
Phone-number
address
email

- b) - Roll-up on Time from day to year  
 - slice for Year 2020  
 - Roll-up on Patient from individual Patient to all

2)  $\text{min\_sup} = 0.6 \times 5 = 3$

a) Apriori Alg.

L<sub>1</sub>

C <sub>1</sub>	itemset	sup_count
for count of each candidate	{M}	3
→	{O}	3
Scan D	{N}	2
for count of each candidate	{K}	5
→	{E}	4
	{Y}	3
	{D}	1
	{A}	1
	{uf}	1

Compare candidate support\_count with minimum support\_count

L <sub>1</sub>	itemset
	{K}
	{E}
	{O}
	{M}
	{Y}

C<sub>2</sub>

Generate C<sub>2</sub>  
candidates

From L<sub>1</sub>

Scan D

For count of  
each candidate

itemset	sup_count
{K, E}	4
{K, O}	3
{K, M}	3
{K, Y}	3
{E, O}	3
{E, M}	2
{E, Y}	2
{O, M}	1
{O, Y}	2
{M, Y}	2

Compare candidate  
support count with  
min\_support count

L<sub>2</sub>

itemset	sup_count
{K, E}	4
{K, O}	3
{K, M}	3
{K, Y}	3
{E, O}	3

C<sub>3</sub>

Generate C<sub>3</sub>  
candidates

From L<sub>2</sub>

Scan D for  
count of each  
candidate

itemset	sup_count
{K, E, O}	3

Compare candidate  
support count with  
min\_support count

L<sub>3</sub>

itemset	sup_count
{K, E, O}	3

$$C_4 = \{\} \Rightarrow L_4 = \{\}$$

Final result:

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

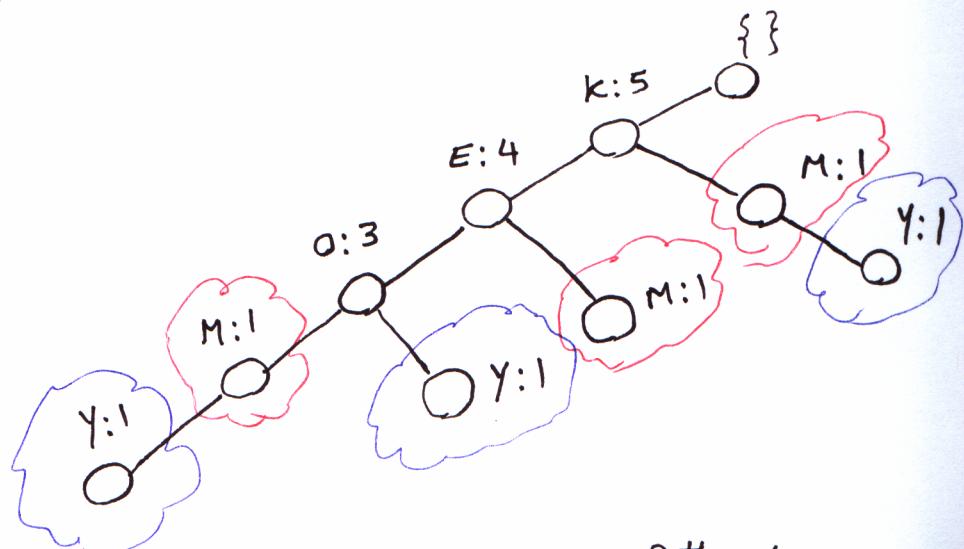
2) FP-tree first scan database & find frequent 1 itemset.

$$L = \{ \{K:5\}, \{E:4\}, \{O:3\}, \{M:3\}, \{Y:3\} \}$$

Then sort frequent items in frequency descending order

TID	items_bought
T100	{K, E, O, M, Y}
T200	{K, E, O, Y}
T300	{K, E, M}
T400	{K, M, Y}
T500	{K, E, O, O}

scan database again &  
construct FP-tree



item	Conditional Pattern Base (CPB)	Conditional FP-tree min_sup = 3	Frequent Patterns Generated (FPS)
Y	{K, E, O, M : 1} {K, E, O : 1} {K, M : 1}	{K:3}	{K, Y: 3}
M	{K, E, O : 1} {K, E : 1} {K : 1}	{K:3}	{K, M : 3}
O	{K, E : 3}	{K, E : 3}	{K, E, O : 3} {E, O : 3} {K, O : 3}
E	{K : 4}	{K:4}	{K, E : 4}

FP-tree result:

$$\left\{ \begin{array}{l} \{K:5\}, \{E:4\}, \{O:4\}, \{M:3\}, \{Y:3\}, \{K,Y:3\}, \\ \{K,M:3\}, \{K,O:3\}, \{K,E:4\}, \{E,O:3\}, \{K,E,O:3\} \end{array} \right\}$$

FP-growth is more efficient because it reduces the size of data set to be searched.

3) a)

$$\text{support} = \frac{2000}{5000} = 0.4 = 40\% > 25\%.$$

$$\text{confidence} = \frac{2000}{3000} = 0.66 = 66.7\% > 50\%.$$

Therefore association rule is strong

b)

$$\text{Lift} = \frac{P\{\text{hot dogs, hamburgers}\}}{P\{\text{hot dogs}\} P\{\text{hamburgers}\}} = \frac{0.4}{0.6 \times 0.5} = 1.3 > 1$$

They are positively correlated.

$$P(\{\text{hot dogs}\}) = \frac{3000}{5000} = 0.6$$

$$P(\{\text{hamburgers}\}) = \frac{2500}{5000} = 0.5$$

$$P(\{\text{hot dogs, hamburgers}\}) = \frac{2000}{5000} = 0.4$$

	hot dogs	hot dogs	$\sum \text{Row}$
hamburgers	2000 (1500)	500 (1000)	2500
hamburgers	1000 (1500)	1500 (1000)	2500
$\sum \text{col}$	3000	2000	5000

2 x 2 table

$$\text{Freedom} = (2-1)(2-1)$$

$$\text{Freedom} = 1$$

$$c_{11} = \frac{\text{count (hot dogs)} \times \text{count (hamburgers)}}{n} = \frac{3000 \times 2500}{5000} = 1500$$

$$c_{12} = \frac{\text{count (hot dogs)} \times \text{count (hamburgers)}}{n} = \frac{2000 \times 2500}{5000} = 1000$$

$$c_{21} = \frac{\text{count (hot dogs)} \times \text{count (hamburgers)}}{n} = \frac{3000 \times 2500}{5000} = 1500$$

$$c_{22} = \frac{\text{count (hot dogs)} \times \text{count (hamburgers)}}{n} = \frac{2000 \times 2500}{5000} = 1000$$

$$\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}} =$$

$$\frac{(2000 - 1500)^2}{1500} + \frac{(500 - 1000)^2}{1000} + \frac{(1000 - 1500)^2}{1500} + \frac{(1500 - 1000)^2}{1000} = 833.33$$

$\chi^2 > 10.828 \Rightarrow$  They are strongly correlated.