

Assignment #3

- 1) Suppose that a data warehouse consists of 3D: time, doctor, patient w/ 2 measures: count, charge where:
- charge = fee doctor charges a patient for a visit

a) Draw a star schema

Visits

Fact table

- time key
- doctors key
- patients key

- charge
- count

Time

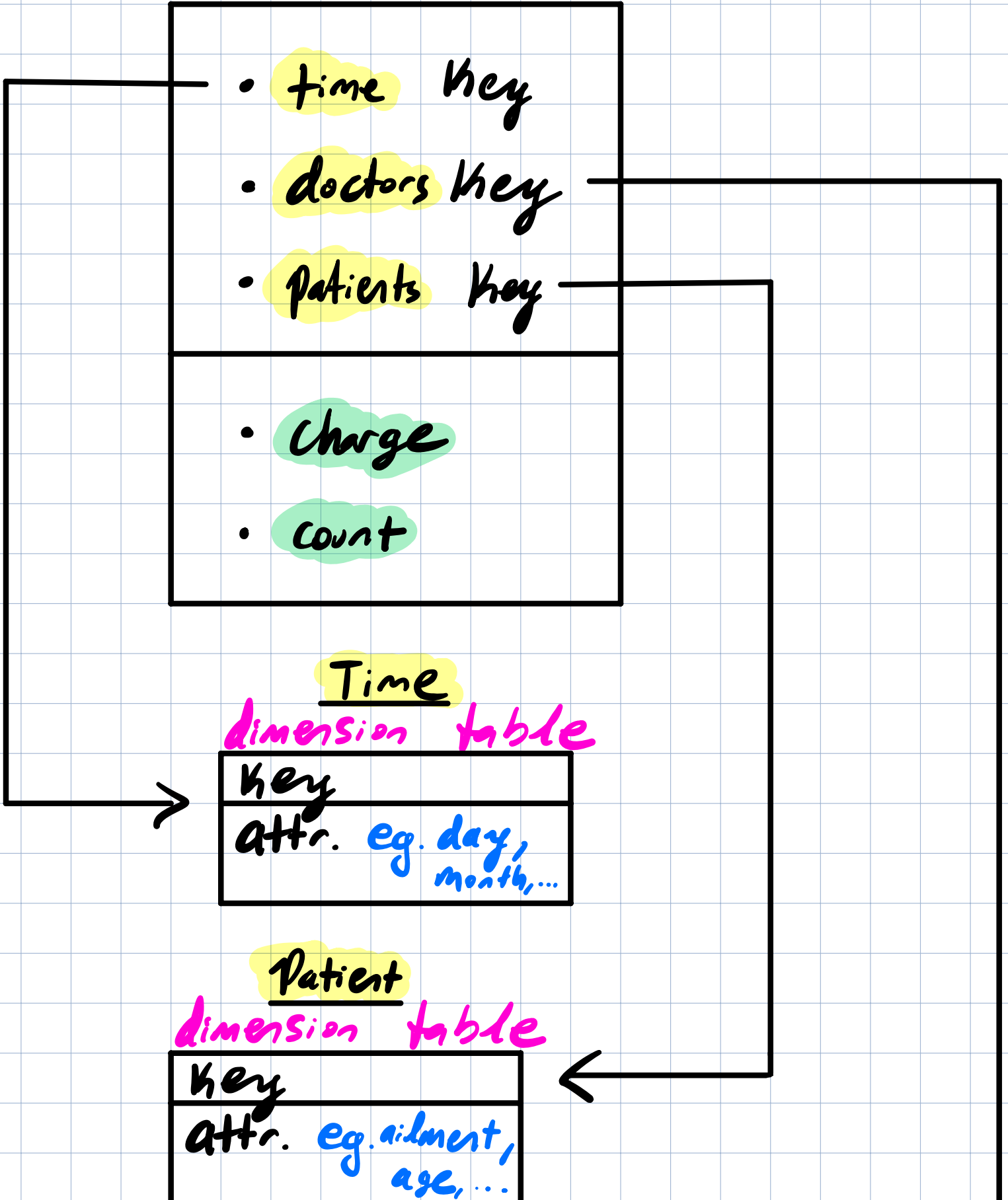
dimension table

key
attr. eg. day, month,...

Patient

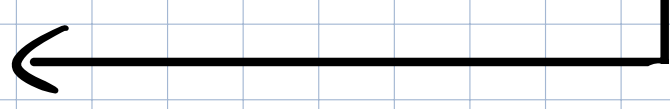
dimension table

key
attr. eg. ailment, age,...




doctor
dimension table

key
attr. eg. specialty, age, ...



- The measures belong to the main fact table.
- The main fact table holds the keys to the dims provided.
- b/c the measure given describes a visit & relies on the data from the dims, we name the fact table "Visits".

b) Starting w/ a base cuboid
[day, doctor, patient], what
time  OLAP ops should be
performed in order to list the
total fee collected \forall doctors
in 2020?

1) Take day & perform a
roll-up on day \rightarrow year.

2) Our data cube now holds
info : [year, doctor, patient].
We Pivot so the data cube

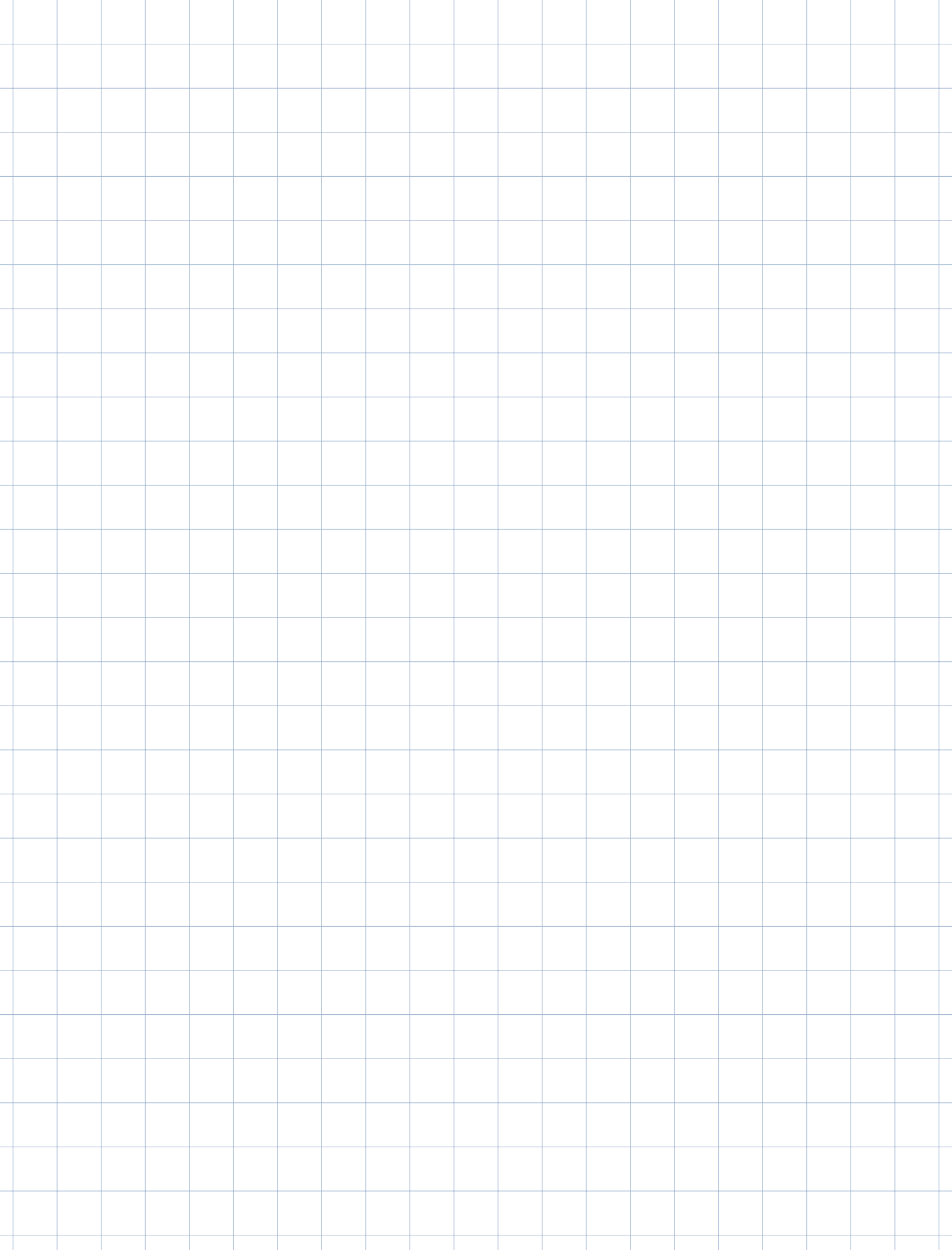
is shaped [doctor, year, patient]

3) W/ this new view, we are able to **slice & Dice** the time/year dimension to just 2020.

4) The resulting data cube has the form [doctor, 2020, patient].

We now use the measure **charge** on this data cube to see the fees collected \forall doctors in 2020.

5) Perform a **sum** operation ^{& OLAP op} to aggregate the charges.



2) A db has 5 transactions

$\text{min_sup} = 60\%$

$\text{min_conf} = 80\%$

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, I, E}

Find all frequent itemsets using

Apriori & Fp-Growth.

Compare efficiency.

$$|D| = 5.$$

$$\text{min_sup} = 60\% \cdot 5 \text{ transactions}$$

$$= 3 \text{ transactions}$$

a) Apriori Alg

L_1

Item	Supp. Ct.
M	3
O	3
N	2
K	5
E	4
Y	3
D	1
A	1
V	1
C	2
I	1

compare
w/
min_sup



i.e. \forall items

$\text{supp-ct} \geq \text{min-sup}$

L_2

Item	Supp. Ct.
M	3
O	3
K	5
E	4
Y	3



L_2

Item Set	Supp. Ct
MO	1
MK	3
ME	2
MY	2
OK	3
OE	3
OY	2
KE	4
KY	3
EY	2

compare
w/ min
sup



L_2

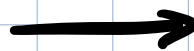
Item Set	Supp Ct.
MK	3
OK	3
OE	3
KE	4
KY	3



L_3

Item Set	Supp. Ct
MOK	1
MOE	1
MKE	2
MKY	2
OKE	3
OKY	2
KEY	2

compare
w/ min
sup



L_3

Item Set	Supp. Ct
OKE	3

b) FP-growth (Tree)

Recall $\text{min-conf} = 80\%$.

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, I, E}

MONKEY

$$M \Rightarrow ONKEY = \frac{1}{3} \times$$

$$ONKEY \Rightarrow M = \frac{1}{2} \times$$

$$O \Rightarrow MNKEY = \frac{1}{3} \times$$

$$MNKEY \Rightarrow O = \frac{1}{1} \checkmark$$

$$N \Rightarrow MOKEY = \frac{1}{2} \times$$

$$MOKEY \Rightarrow N = \frac{1}{1} \checkmark$$

$$K \Rightarrow MONEY = \frac{1}{5} \times$$

$$MONEY \Rightarrow K = \frac{1}{1} \checkmark$$

$$E \Rightarrow MONKY = \frac{1}{4} \times$$

$$MONKY \Rightarrow E = \frac{1}{1} \checkmark$$

$$Y \Rightarrow MONKE = \frac{1}{3} \times$$

$$MONKE \Rightarrow Y = \frac{1}{1} \checkmark$$

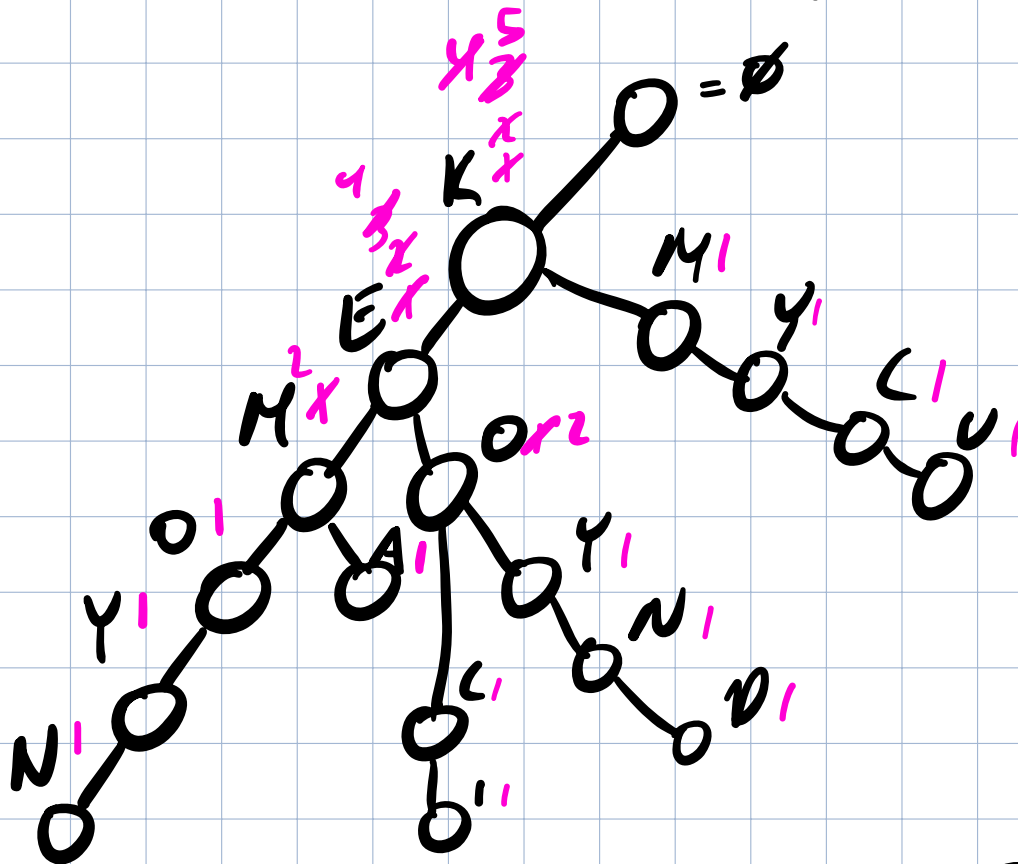
⋮

FP Tree

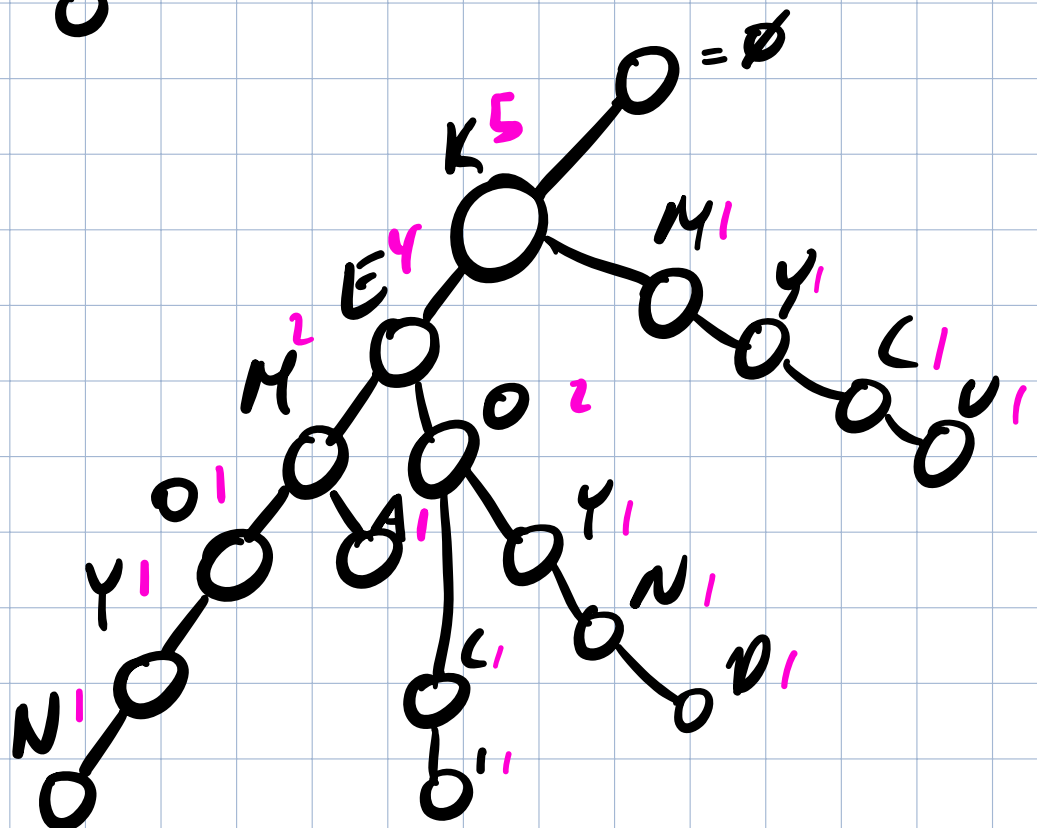
Sort C_1 by Supp-ct \forall items

\Rightarrow KEMOYNLDAVI

~~DONKEY~~
MAKE
~~MUCKY~~
LOOKIE



\Rightarrow



FP Tree is easy to make.

3) Contingency table given

	hot dogs	$\overline{\text{hot dogs}}$	Σ_{row}
hamburgers	2000	500	2500
$\overline{\text{hamburgers}}$	1000	1500	2500
Σ_{col}	3000	2000	5000

item $\in \{\text{hot dogs, hamb}\}$

item refers to transactions containing item

a) Assoc Rule: hot dogs \Rightarrow hamburgers.

min-sup = 25%, min-conf = 50%

Strong?

✓ $\frac{2000}{5000} = \frac{2}{5} \geq \frac{1}{4} = 25\% = \text{min-sup}$

✓ $\frac{2000}{3000} = \frac{2}{3} \geq \frac{1}{2} = 50\% = \text{min-conf}$

\therefore The given Assoc Rule is strong
w/ the provided min-sup & min-conf.

b) Calculate lift & Chi-square.

$$\text{The call lift} = \frac{P(A \cup B)}{P(A) \cdot P(B)}$$

$$P(\text{hot dogs}) = \frac{3000}{5000} = \frac{3}{5}$$

$$P(\text{hamburgers}) = \frac{2500}{5000} = \frac{1}{2}$$

$$P(\text{hot dogs, hamburgers}) = \frac{2000}{5000} = \frac{2}{5}$$

$$\text{Lift} = \frac{\left(\frac{2}{5}\right)}{\left(\frac{3}{5}\right) \cdot \left(\frac{1}{2}\right)} = \frac{\left(\frac{2}{5}\right)}{\left(\frac{3}{10}\right)} = \left(\frac{2}{\cancel{5}}\right) \cdot \left(\frac{10^2}{3}\right) = \frac{4}{3} > 1$$

\therefore hot dogs & hamburgers
are positively corr.

the call χ^2 formula

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

	hot dogs	$\overline{\text{hot dogs}}$	Σ_{row}
hamburgers	2000	500	2500
$\overline{\text{hamburgers}}$	1000	1500	2500
Σ_{col}	3000	2000	5000

$$C_{11} = \frac{3000 \cdot 2500}{5000} = 1500$$

$$C_{12} = \frac{2000 \cdot 2500}{5000} = 1000$$

$$C_{21} = \frac{3000 \cdot 2500}{5000} = 1500$$

$$C_{22} = \frac{2000 \cdot 2500}{5000} = 1000$$

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

$$\Rightarrow \chi^2 = \frac{2500}{3} = 833.\overline{3}$$

Recall Degree of Freedom = $(dim 1 - 1) \cdot (dim 2 - 1)$

$$\Rightarrow DF = 1.$$

max prob. is 10.828

\Rightarrow we reject the hypothesis
that hot dogs & hamburgers
are independent.

\Rightarrow they're dependent.

c) What kind of corr. relationship
exists b/w purchase of hot dogs &
purchase of hamburgers.

They're positively corr. as found
by lift & verified by x^2
in part b).