

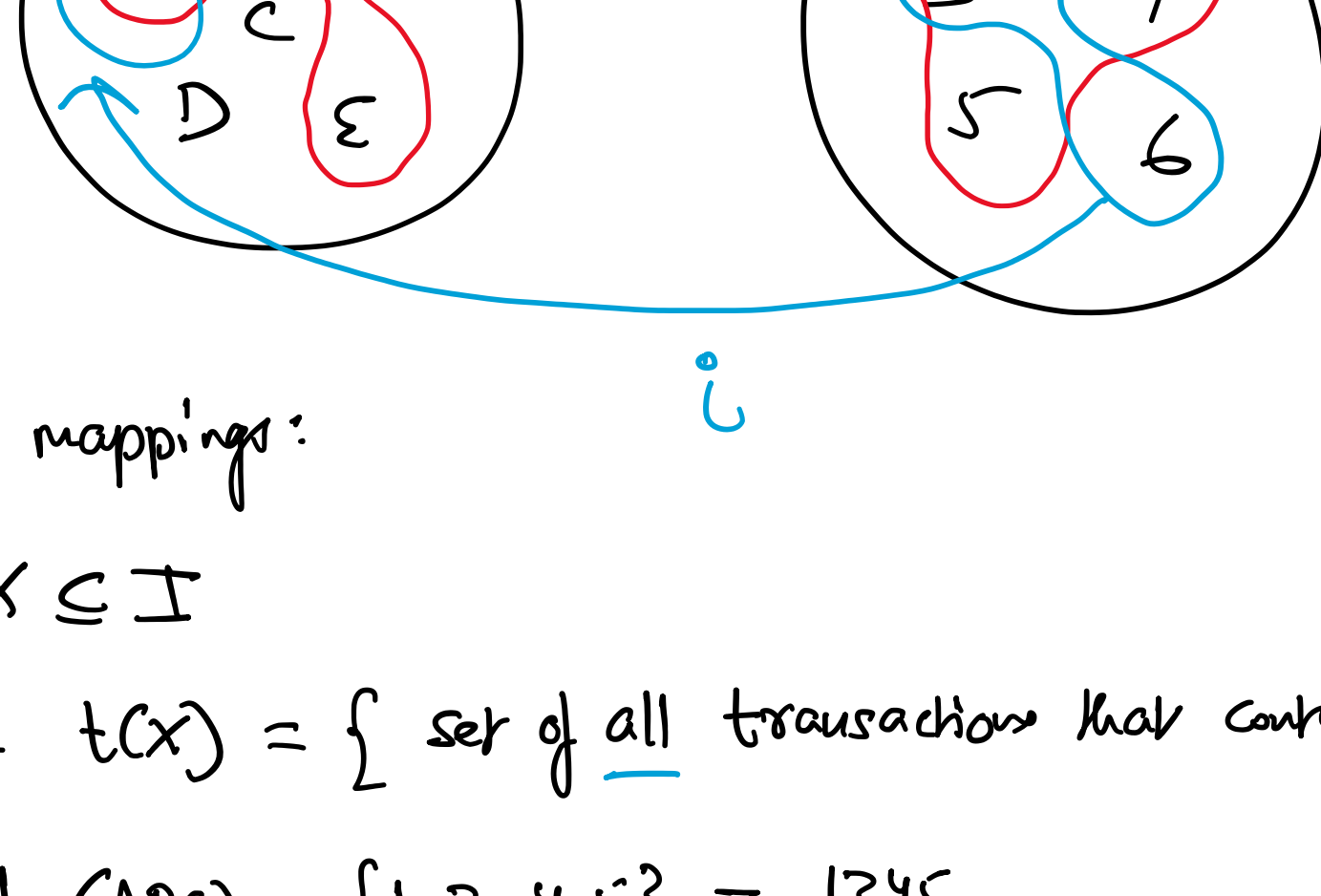
frequent itemsets

$$X \subseteq I$$

all the items

$$Y \subseteq T$$

set of samples / instance / transactions ID



- dataset
- ABDE
 - BCE
 - ABDE
 - ABCE
 - ABCEDE
 - BCE

two mappings:

$$X \subseteq I$$

define $t(X) = \{ \text{set of all transactions that contain } X \}$

$$t(\{A, B, E\}) = t(ABE) = \{1, 3, 4, 5\} = 1345$$

$t(X) \equiv$ tidset for X

$$\text{sup}(X) = |t(X)|$$

$\uparrow \equiv$ # of transactions containing X

support of X

frequency of X

$$\hat{p}(X) = \frac{\text{sup}(X)}{|T|}$$

estimate for prob of X
joint probability of all the items in X

Counting

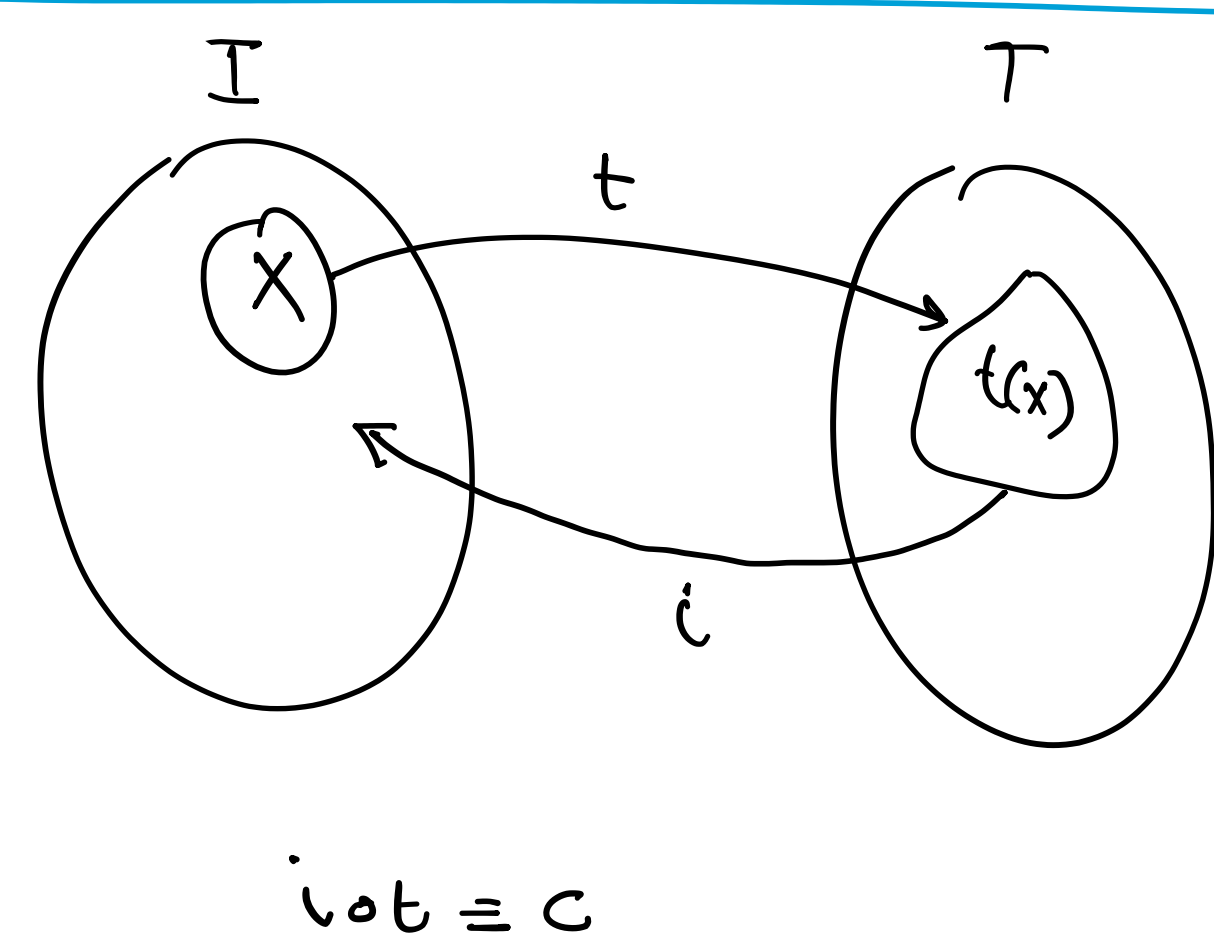
$$\hat{p}(ABE) = \frac{4}{6} = \frac{2}{3} = 0.67$$

observed prob

expected prob under independence

$$\hat{p}(ABE) \neq \hat{p}(A) \times \hat{p}(B) \times \hat{p}(E)$$

- ABDE
- BCE
- ABDE
- ABCE
- ABCEDE
- BCE



$$t(ABE) = 1345$$

$$i(t(ABE)) = i(1345) = ABE$$

$$i(t(ABE)) = i(1345) = ABE$$

$$i \circ t \equiv C$$

Closure operation

If $C(X) = X$ then we say that X is a closed set

$i \circ t(X)$

Comparison

e.g. ABE is a closed set

$X = ABC$ not closed since $C(ABC) = ABCE$

$C(X)$ is always closed, (even if X is not closed)

Closure operator C :

1) idempotent $C(X) = C(C(X))$

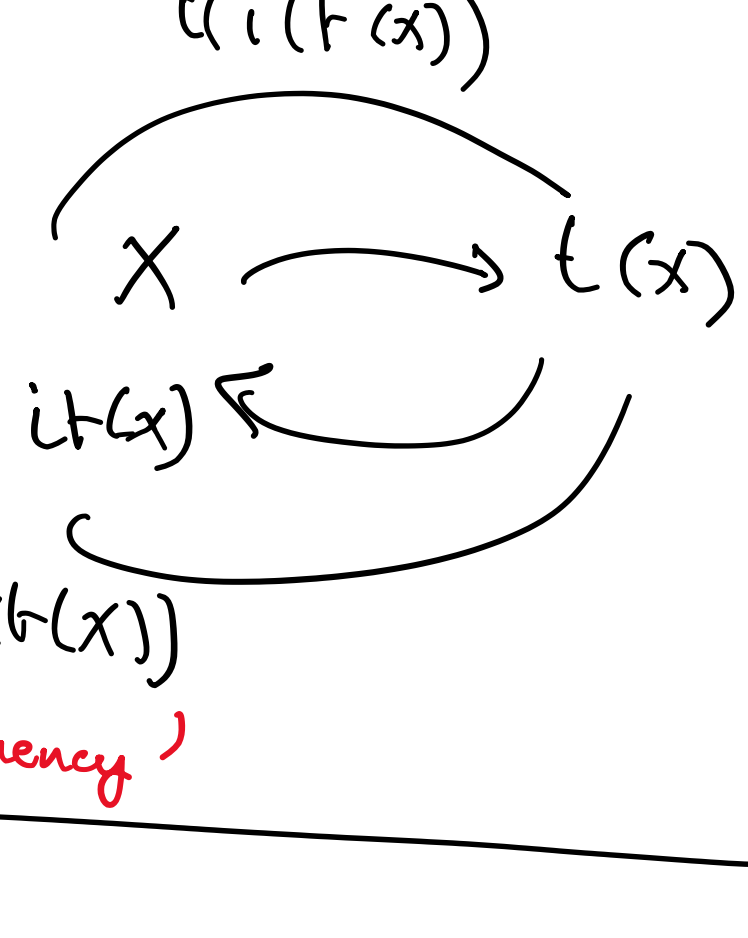
2) extensive $C(X) \supseteq X$

3) monotone

$$X_1 \subseteq X_2$$

$$\Rightarrow C(X_1) \subseteq C(X_2)$$

X is closed if there is no superset with the same frequency



Exclar: intersections of tidsets

- ABDE
- BCE
- ABDE
- ABCE
- ABCEDE
- BCE

$$\Rightarrow$$

$$x \rightarrow t(x)$$

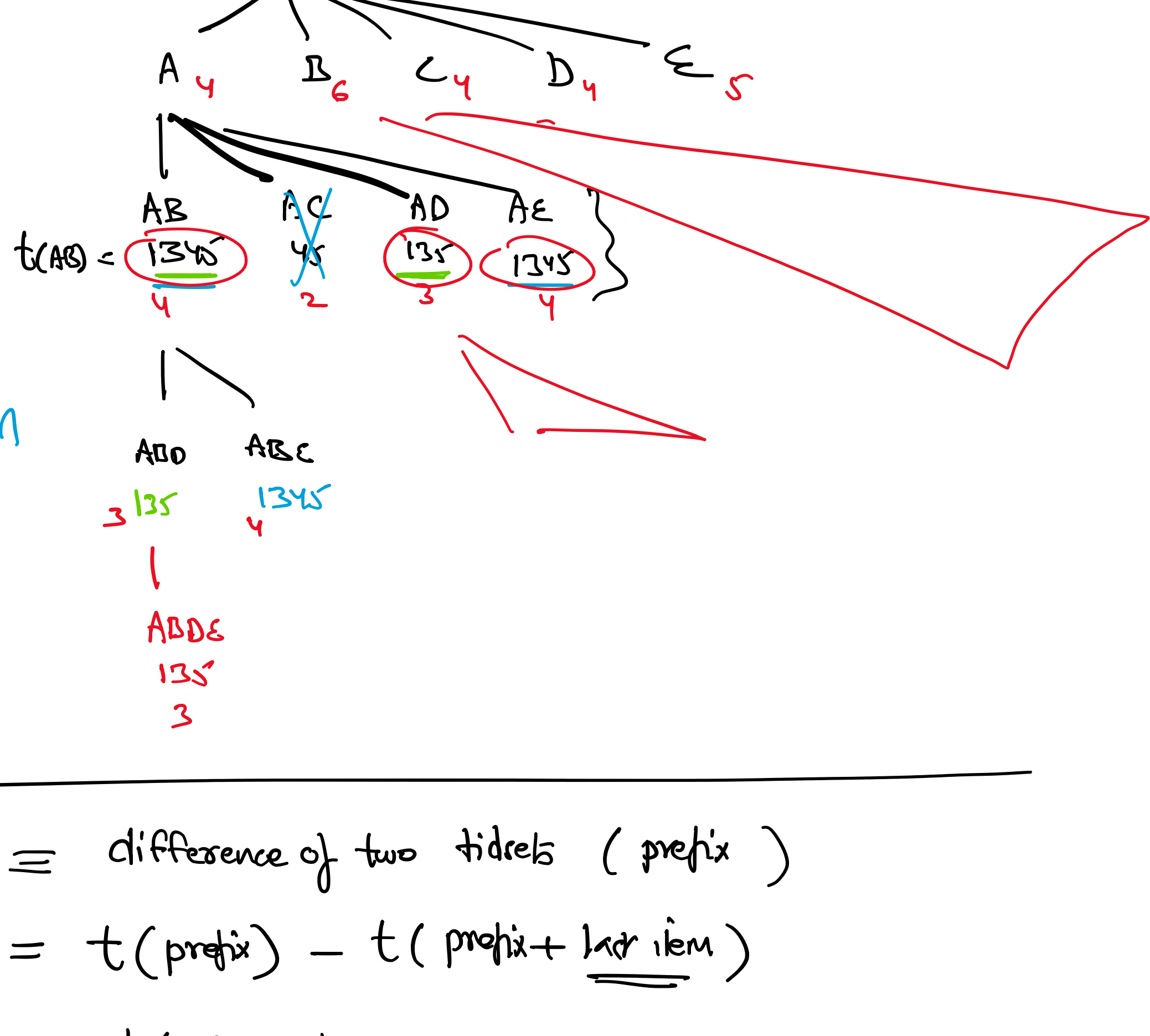
x	t(x)
A	1345
B	123456
C	2456
D	1356
E	12345

Inverted dataset

Combinatorial Search

$$t(AB) = t(A) \cap t(B)$$

Minsup = 3



$$t(ABD) = t(AB) \cap t(AD)$$

diffset \equiv difference of two tidsets (prefix)

$$= t(\text{prefix}) - t(\text{prefix} + \text{last item})$$

$$d(ABD) = t(AB) - t(ABD)$$

$$\text{initially: } d(A) = t(\emptyset) - t(A)$$

$$A = \{\emptyset, A\}$$

$$= 123456 - 1345$$

$$= 26$$

those transactions that do not contain A

$$\text{sup}(B) = \text{sup}(\emptyset) - \text{len}(d(A))$$

all tids

$$123456 - 26$$

$$1345$$

all tids that contain A

$$6 - 2 = 4 \text{ sup}$$

$$d(X) = t(\text{prefix}) - t(X)$$

diffsets from diffset

$$d(\text{prefix} + \text{item1} + \text{item2}) = d(\text{prefix} + \text{item2}) - d(\text{prefix} + \text{item1})$$

$$d(ABD) = d(AD) - d(AB)$$

||

$$\emptyset - t(\emptyset) = 123456$$

$$t(AB) - t(ABD)$$

declar

X	t(X)	d(X)
A	1345	26
B	123456	6
C	2456	13
D	1356	24
E	12345	6

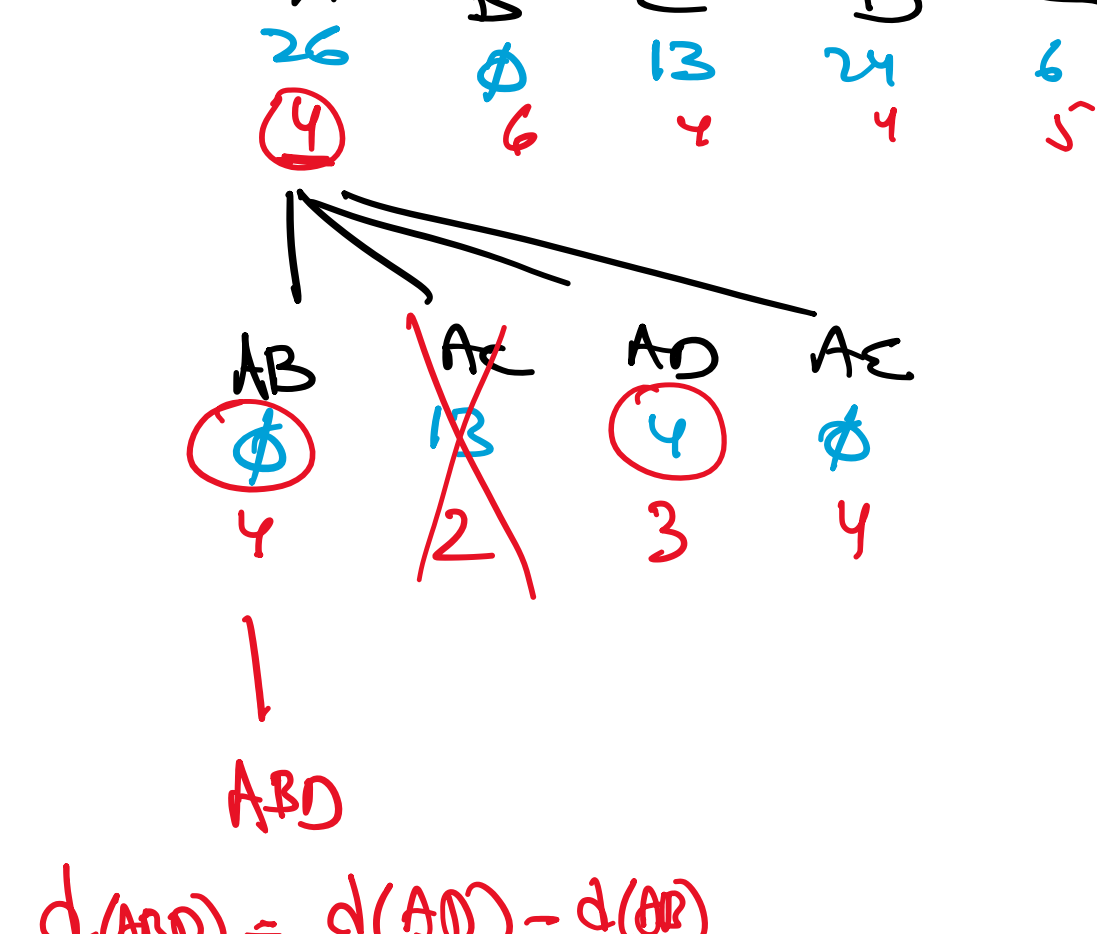
$$d(AD) = d(B) - d(AB)$$

$$= 6 - \{2, 6\}$$

$$= \emptyset$$

$$\{1, 3\} - \{2, 6\}$$

$$= \{1, 3\}$$



$$d(ABD) = d(AD) - d(AB)$$

$$= 4 - \emptyset$$

$$= 4$$

$$\text{sup} = 3$$

closed sets are loss-less compressed representation of all freq sets

F = set of all frequent itemsets

- ABDE
- BCE
- ABDE
- ABCE
- ABCEDE
- BCE

Minsup = 3

No superset with same freq \Rightarrow closed

freq	itemsets
6	B
5	E, BE
4	A, C, D, AB, AE, AD, BD, ABDE
3	AD, CE, DE, ABD, ADE, BCE, BDE, ABDE

F

19 frequent itemsets

$$C(B) = B$$

\mathcal{C} = Set of all closed and frequent itemsets

$$= \{B, BE, BC, BD, AB, AE, AD, BD, BCE, BDE, ABDE\}$$

6

5

4

4

4

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

3

</