

# Frequent Itemset

8/31

Let  $I = \{x_1, \dots, x_n\}$  items (products in the grocery)

$X \subseteq I$  is called an item set

itemset of size  $k$   $k$ -itemset

{milk, eggs}

↑  
2-itemset

$I^{(k)}$  set of all  $k$ -itemsets

$T = \{t_1, \dots, t_n\}$  a set of transactionids tids

txn

$T \subseteq T$  is called tidset  
(assumed sorted lex order)

Let  $t$  be a unique txn identifier  
 $X$  be an itemset

A transaction  $(t, X)$

A set of txs

{

(abc123, {milk, eggs})  
(abdl, {carrots, peaches, bananas})

,

,

,

}

## Database Rep

binary DB  $D \subseteq T \times L$

tid  $t \in T$   
item  $x \in L$

$(t, x) \in D \Leftrightarrow$  for some trans  $(t, x)$   $x \in X$

$t$  contains item  $x \Leftrightarrow (t, x) \in D$

$$\text{E.g. } L = \{A, B, C, D, E\} \\ T = \{1, 2, 3, 4, 5, 6\}$$

$\begin{pmatrix} 1, ABDE \\ 2, BCE \end{pmatrix}$

	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1

For a set  $X$  let  $2^X$  be the power set of  $X$

Let  $i: 2^T \rightarrow 2^X$

$$T \subseteq T$$

$$i(T) = \{x \mid \forall t \in T, t \text{ contains } x\}$$

all items in all trans in  $T$

$$\text{E.g. } i(\{2, 3, 4\}) = \{B, E\}$$

Sometimes write transDB as tuples  $(t, i(t))$

D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

# Running Example

binary data base

<b>D</b>	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

<b>t</b>	<b>i(t)</b>
1	<i>ABDE</i>
2	<i>BCE</i>
3	<i>ABDE</i>
4	<i>ABCE</i>
5	<i>ABCDE</i>
6	<i>BCD</i>

let  $t: 2^{\mathbb{F}} \rightarrow 2^{\mathbb{F}}$

$$X \subseteq \mathbb{F}$$

$$t(X) = \{t \mid t \in T \text{ and } t \text{ contains } X\}$$

all of the tids containing all items of itemset

$$t(\{A, C, E\}) = \{4, 5\}$$

<b>D</b>	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

for item  $x \in I$

for all tids

w/ item  $x$

$$t(x)$$

# Vertical DB

for item  $x \in \Sigma$   
tuples  $(x, t(x))$

## Binary DB

D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

## Vertical DB

x	A	B	C	D	E
t(x)	1	1	2	1	1
	3	2	4	3	2
	4	3	5	5	3
	5	4	6	6	4
		5			5
		6			

Support of itemset  $X$  in  $D$   $\sup(X, D)$   $\rightarrow \sup(X)$   
 when database  
 is clear from  
 context

$$\sup(X, D) = |\{t \mid (t, i(t)) \in D \text{ and } X \subseteq i(t)\}| = |t(X)|$$

relative support of  $X$

$$rsup(X, D) = \frac{\sup(X, D)}{|D|}$$

$\downarrow rsup(X)$  when  $D$  is clear from context  
 (est of joint prob of items containing  $X$ )

given user defined min support threshold  $(\text{minsup})$

itemset  $X$   
 frequent in  $D$

if  $\sup(X, D) \geq \text{minsup}$  if  $\text{minsup} \in \mathbb{Z}^+$

$rsup(X, D) \geq \text{minsup}$  if  $\text{minsup} \in [0, 1]$

$F$  is the set of all frequent itemsets

$F^{(k)}$  is the set of all frequent k-itemsets

D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

t	i(t)
1	ABDE
2	BCE
3	ABDE
4	ABCE
5	ABCDE
6	BCD

$$\text{minSup} = 3$$

$$F^{(1)} = ?$$

$$F^{(2)} = ?$$

$$F^{(3)} = ?$$

↑ think about

$$\text{sup}(ABCDE)$$