

# IS422P - DATA MINING

## 6- MINING FREQUENT PATTERNS, ASSOCIATIONS, & CORRELATIONS



AMIRA REZK  
INFORMATION SYSTEM DEPARTMENT



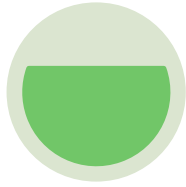


## The Basics

Market Basket  
Analysis

Frequent Itemsets

Association Rules

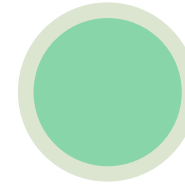


## Frequent Itemset Mining Methods

Apriori Algorithm

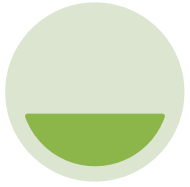
Generating  
Association Rules  
from Frequent  
Itemsets

FP-Growth



## Pattern Evaluation Methods



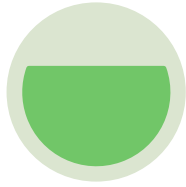


## The Basics

Market Basket  
Analysis

Frequent Itemsets

Association Rules

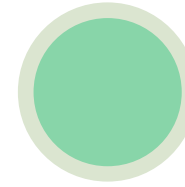


Frequent Itemset  
Mining Methods

Apriori Algorithm

Generating  
Association Rules  
from Frequent  
Itemsets

FP-Growth



Pattern Evaluation  
Methods



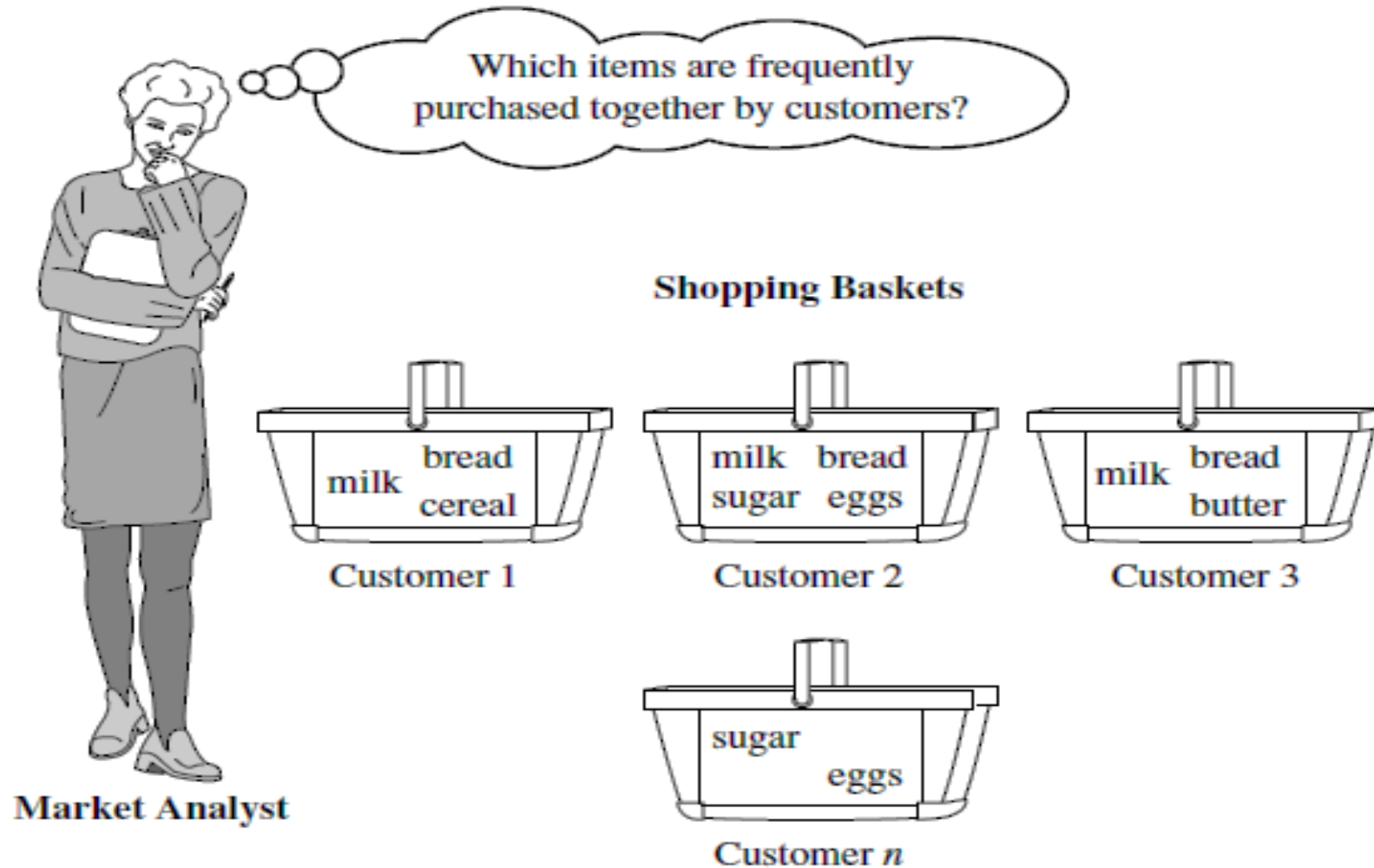
## THE BASICS

### WHAT IS FREQUENT PATTERN ANALYSIS?

- Frequent pattern: a pattern (a set of items, subsequences, substructures, etc.) that occurs frequently in a data set
- First proposed by Agrawal, Imielinski, and Swami [AIS93] in the context of frequent itemsets and association rule mining



# THE BASICS



# THE BASICS

- **Motivation:** Finding inherent regularities in data
  - What products were often purchased together?— Beer and diapers?!
  - What are the subsequent purchases after buying a PC?
  - What kinds of DNA are sensitive to this new drug?
  - Can we automatically classify web documents?
- **Applications**
  - Basket data analysis, cross-marketing, catalog design, sale campaign analysis, Web log (click stream) analysis, and DNA sequence analysis



# THE BASICS

- Frequent Pattern are itemsets that appear frequently in a data set (e.g.Transaction record)
- Items that are frequently associated (e.g purchased) together can be represented as association rules

**Computer → antivirus\_SW [Support = 2% , Confidence =60%]**

- **Support** and **Confidence** are measures of rule interestingness
- 2% Support means 2% of Transactions Show that computers and antivirus\_SW are bought Together
- 60% Confidence means 60 % of customers who bought a computer also bought antivirus\_SW



# THE BASICS

## FREQUENT ITEM-SETS

- Itemset  $X = \{x_1, \dots, x_k\}$                       ex:  $X = \{A, B, C, D, E, F\}$
- Find all the rules  $X \rightarrow Y$  with minimum support and confidence
  - support,  $s$ , probability that a transaction contains  $X \cup Y$
  - confidence,  $c$ , conditional probability that a transaction having  $X$  also contains  $Y$

$$\text{support } X \rightarrow Y = P(X \cup Y) = \frac{n(X \cup Y)}{N}$$

$$\text{confidence } (X \rightarrow Y) = P(Y|X) = \frac{n(X \cup Y)}{n(X)}$$





# THE BASICS

## ASSOCIATION RULES

Ex: Let  $\text{min\_Sup.} = 50\%$ ,  $\text{min\_conf.} = 50\%$

*Frequent Patterns:*

$\{A:3, B:3, D:4, E:3, AD:3\}$

Association rules:

$A \rightarrow D$  (60%, 100%)

$D \rightarrow A$  (60%, 75%)

$$\text{conf}(A \rightarrow D) = \frac{3}{3} = 100\%$$

$$\text{conf}(D \rightarrow A) = \frac{3}{4} = 75\%$$

Transaction-id	Items bought
10	A, B, D
20	A, C, D
30	A, D, E
40	B, E, F
50	B, C, D, E, F



# THE BASICS

## ASSOCIATION RULES

- If frequency of itemset  $I$  satisfies  $\text{min\_support}$  count then  $I$  is a frequent itemset
- If a rule satisfies  $\text{min\_support}$  and  $\text{min\_confidence}$  thresholds, it is said to be strong
  - problem of mining association rules reduced to mining frequent itemsets
- Association rules mining becomes a two-step process:
  - Find all frequent itemsets that occur at least as frequently as a predetermined  $\text{min\_support}$  count
  - Generate strong association rules from the frequent itemsets that satisfy  $\text{min\_support}$  and  $\text{min\_confidence}$



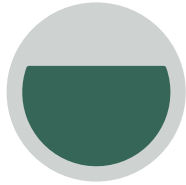


## The Basics

Market Basket  
Analysis

Frequent Itemsets

Association Rules

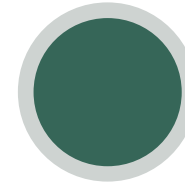


## Frequent Itemset Mining Methods

Apriori Algorithm

Generating  
Association Rules  
from Frequent  
Itemsets

FP-Growth



## Pattern Evaluation Methods



# MINING FREQUENT ITEMSETS

## APRIORI ALGORITHM

- Goes as follows:
  - Find frequent 1-itemsets  $\rightarrow L_1$
  - Use  $L_1$  to find frequent 2-itemsets  $\rightarrow L_2$
  - ... until no more frequent k-itemsets can be found
- Each  $L_k$  itemset requires a full dataset scan
- To improve efficiency, use the Apriori property:
  - “All nonempty subsets of a frequent itemset must also be frequent” – if a set cannot pass a test, all of its supersets will fail the same test as well – if  $P(I) < \text{min\_support}$  then  $P(I \cup A) < \text{min\_support}$



# MINING FREQUENT ITEMSETS

## APRIORI ALGORITHM

Transactional data example  
N=10, min\_supp count=2

TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2

Scan dataset for  
count of each  
candidate

$C_1$

Itemset	Support count
{I1}	7
{I2}	8
{I3}	6
{I4}	2
{I5}	2

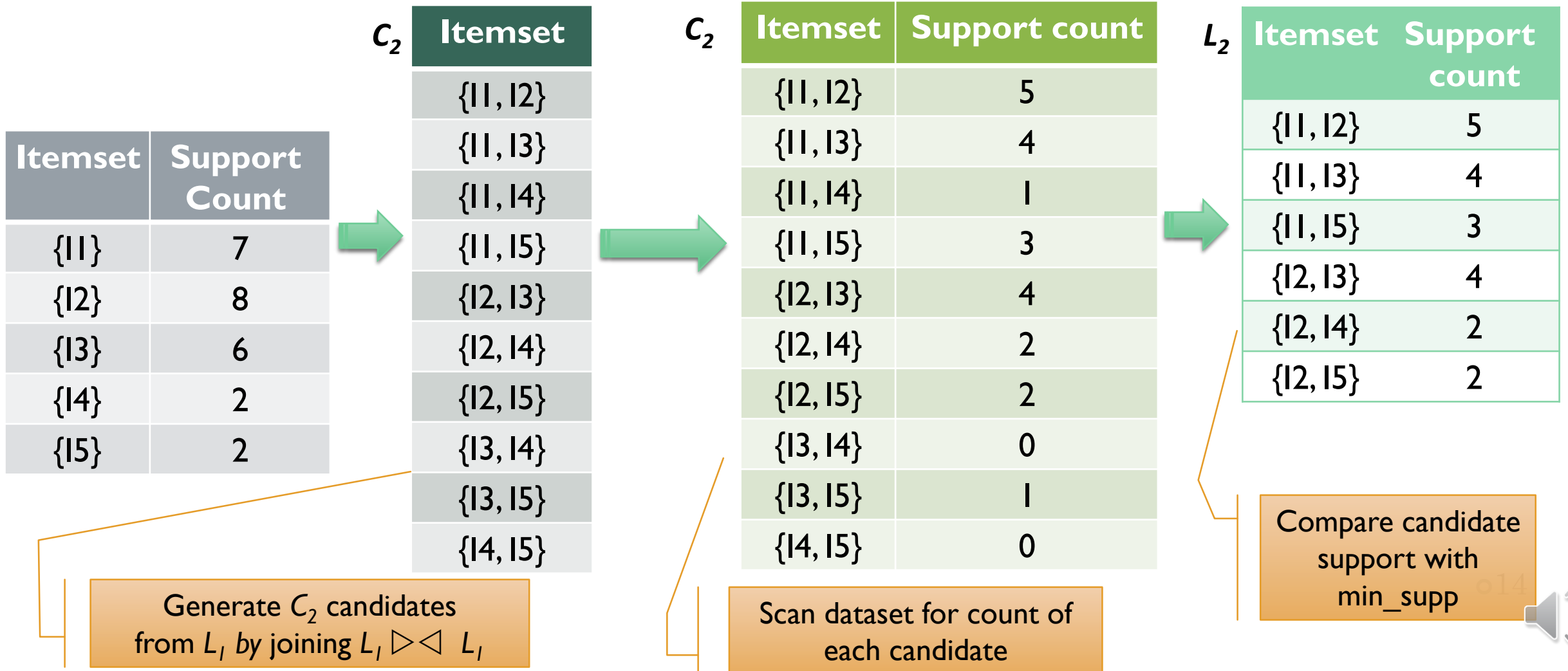
Compare candidate  
support with  
min\_support

$L_1$

Itemset	Support Count
{I1}	7
{I2}	8
{I3}	6
{I4}	2
{I5}	2

# MINING FREQUENT ITEMSETS

## APRIORI ALGORITHM



# MINING FREQUENT ITEMSETS

## APRIORI ALGORITHM

$C_3 = L_2 \triangleright \triangleleft L_2 = \{ \{I1, I2, I3\}, \{I1, I2, I5\}, \{I1, I3, I5\}, \{I2, I3, I4\}, \{I2, I3, I5\}, \{I2, I4, I5\} \}$

Not all subsets are frequent  $\rightarrow$  **Prune** (Apriori property)

Itemset	Support count
{I1, I2}	5
{I1, I3}	4
{I1, I5}	3
{I2, I3}	4
{I2, I4}	2
{I2, I5}	2



Itemset
{I1, I2, I3}
{I1, I2, I5}



$C_3$	
Itemset	Support count
{I1, I2, I3}	2
{I1, I2, I5}	2

Scan dataset for count of each candidate



$L_3$	
Itemset	Support count
{I1, I2, I3}	2
{I1, I2, I5}	2

Compare candidate support with min\_supp

Generate  $C_3$  candidates from  $L_2$  by **joining**  $L_2 \triangleright \triangleleft L_2$



Two joining (lexicographically ordered) k-itemsets must share first k-1 items  $\rightarrow \{I1, I2\}$  is not joined with  $\{I2, I4\}$



# MINING FREQUENT ITEMSETS

## APRIORI ALGORITHM

Itemset	Support count
{11, 12, 13}	2
{11, 12, 15}	2



Itemset
{11, 12, 13, 15}



Not all subsets are frequent → **Pruning**

$C_4 = \phi \rightarrow$  **Terminate**





# APRIORI ALGORITHM

**Algorithm: Apriori.** Find frequent itemsets using an iterative level-wise approach based on candidate generation.

**Input:**

- $D$ , a database of transactions;
- $min\_sup$ , the minimum support count threshold.

**Output:**  $L$ , frequent itemsets in  $D$ .

**Method:**

Generate  $C_k$  using  $L_{k-1}$  to find  $L_k$

```
(1)  $L_1 = \text{find\_frequent\_1-itemsets}(D)$ ;  
(2) for ( $k = 2; L_{k-1} \neq \phi; k++$ ) {  
(3)    $C_k = \text{apriori\_gen}(L_{k-1})$ ;  
(4)   for each transaction  $t \in D$  { // scan  $D$  for counts  
(5)      $C_t = \text{subset}(C_k, t)$ ; // get the subsets of  $t$  that are candidates  
(6)     for each candidate  $c \in C_t$   
(7)        $c.\text{count}++$ ;  
(8)   }  
(9)    $L_k = \{c \in C_k | c.\text{count} \geq min\_sup\}$   
(10) }  
(11) return  $L = \cup_k L_k$ ;
```

Join

```
procedure  $\text{apriori\_gen}(L_{k-1}:\text{frequent } (k-1)\text{-itemsets})$   
(1) for each itemset  $l_1 \in L_{k-1}$   
(2)   for each itemset  $l_2 \in L_{k-1}$   
(3)     if ( $(l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2])$   
        $\wedge \dots \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$ ) then {  
(4)        $c = l_1 \bowtie l_2$ ; // join step: generate candidates  
(5)       if  $\text{has\_infrequent\_subset}(c, L_{k-1})$  then  
(6)         delete  $c$ ; // prune step: remove unfruitful candidate  
(7)       else add  $c$  to  $C_k$ ;  
(8)     }  
(9) return  $C_k$ ;
```

Prune

```
procedure  $\text{has\_infrequent\_subset}(c:\text{candidate } k\text{-itemset};$   
    $L_{k-1}:\text{frequent } (k-1)\text{-itemsets})$ ; // use prior knowledge  
(1) for each  $(k-1)$ -subset  $s$  of  $c$   
(2)   if  $s \notin L_{k-1}$  then  
(3)     return TRUE;  
(4) return FALSE;
```



# MINING FREQUENT ITEMSETS

## GENERATING ASSOCIATION RULES FROM FREQUENT ITEMSETS

- Association rules can be generated using the confidence equation, as follows”

$$\text{confidence}(A \Rightarrow B) = P(B|A) = \frac{\text{support\_count}(A \cup B)}{\text{support\_count}(A)}$$

- For each frequent itemset L, generate all nonempty subset of L
- For every nonempty subset S of L, output rule:  $S \rightarrow L-S$
- $\frac{\text{Supportcount}(L)}{\text{Supportcount}(S)} \geq \text{min\_conf}$ ,

where min\_conf is the minimum confidence threshold.



# MINING FREQUENT ITEMSETS

## GENERATING ASSOCIATION RULES FROM FREQUENT ITEMSETS

Itemset	Support count
{11, 12, 13}	2
{11, 12, 15}	2

Nonempty subsets
{11, 12}
{11, 15}
{12, 15}
{11}
{12}
{15}

Association Rules
$\{11, 12\} \rightarrow 15$
$\{11, 15\} \rightarrow 12$
$\{12, 15\} \rightarrow 11$
$11 \rightarrow \{12, 15\}$
$12 \rightarrow \{11, 15\}$
$15 \rightarrow \{11, 12\}$

Confidence
$2/5 = 40\%$
$2/2 = 100\%$
$2/2 = 100\%$
$2/7 = 28\%$
$2/8 = 25\%$
$2/2 = 100\%$

For a *min\_confidence* = 70%



# MINING FREQUENT ITEMSETS

## FP-GROWTH

- To avoid costly candidate generation
- Divide-and-conquer strategy:
- Compress database representing frequent items into a frequent pattern tree (FP-tree) – 2 passes over dataset
- Divide compressed database (FP-tree) into conditional databases, then mine each for frequent itemsets – traverse through the FP-tree



# MINING FREQUENT ITEMSETS

## FP-GROWTH

Transactional data example  
N=10, min\_supp count=2

TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2

Scan dataset for count  
of each candidate

$C_1$

Itemset	Support count
{I1}	7
{I2}	8
{I3}	6
{I4}	2
{I5}	2

Compare candidate  
support with min\_supp

$L_1$  - Reordered

Itemset	Support count
{I2}	8
{I1}	7
{I3}	6
{I4}	2
{I5}	2



# MINING FREQUENT ITEMSETS

## FP-GROWTH – FP-TREE CONSTRUCTION

$L_1$  - Reordered

Itemset	Support count
{I2}	8
{I1}	7
{I3}	6
{I4}	2
{I5}	2

*FP-tree*



TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2

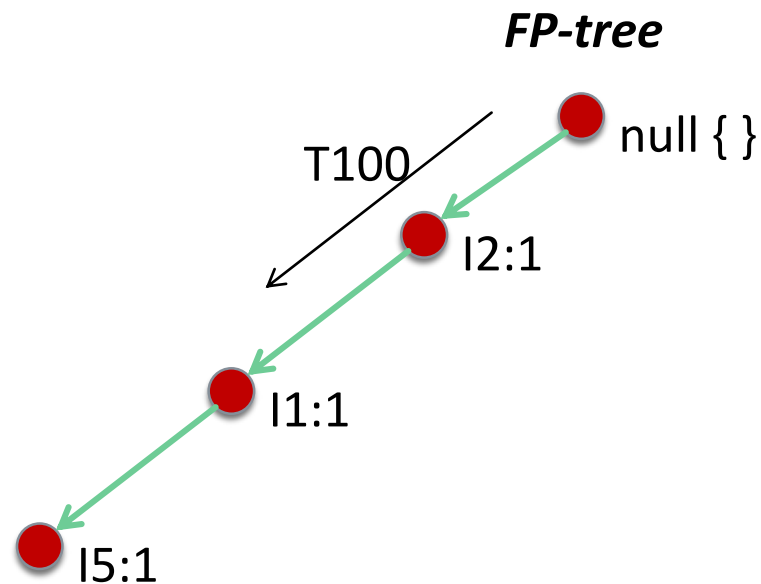


# MINING FREQUENT ITEMSETS

## FP-GROWTH – FP-TREE CONSTRUCTION

$L_1$  - Reordered

Itemset	Support count
{I2}	8
{I1}	7
{I3}	6
{I4}	2
{I5}	2



Order of items is kept throughout path construction, with common prefixes shared whenever applicable

TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2

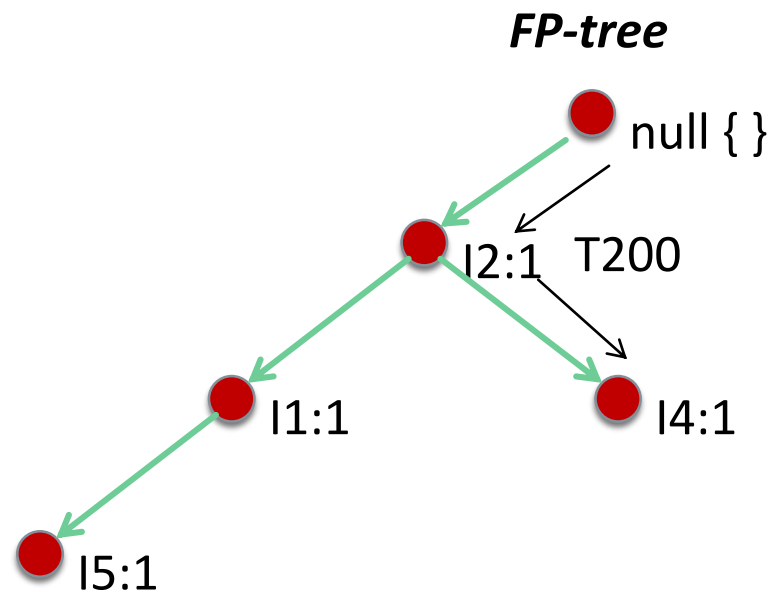


# MINING FREQUENT ITEMSETS

## FP-GROWTH – FP-TREE CONSTRUCTION

$L_1$  - Reordered

Itemset	Support count
{I2}	8
{I1}	7
{I3}	6
{I4}	2
{I5}	2



TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2



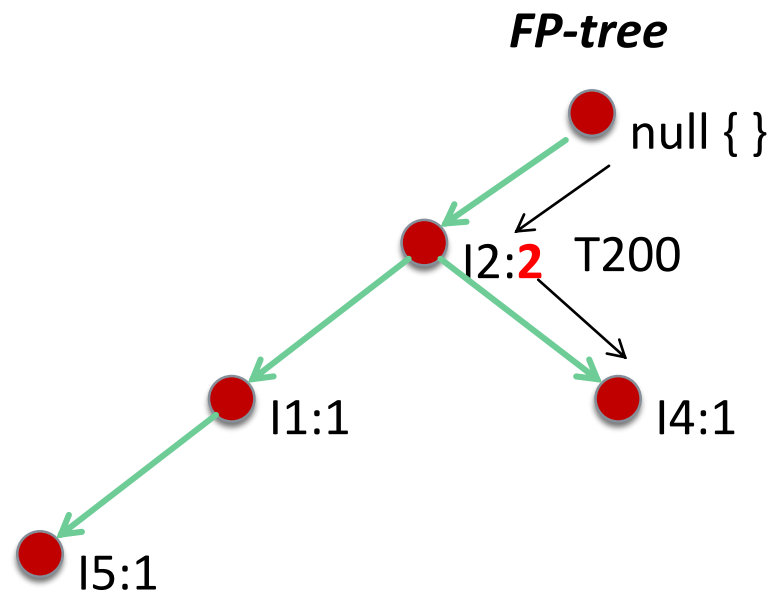


# MINING FREQUENT ITEMSETS

## FP-GROWTH – FP-TREE CONSTRUCTION

$L_1$  - Reordered

Itemset	Support count
{I2}	8
{I1}	7
{I3}	6
{I4}	2
{I5}	2



TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2

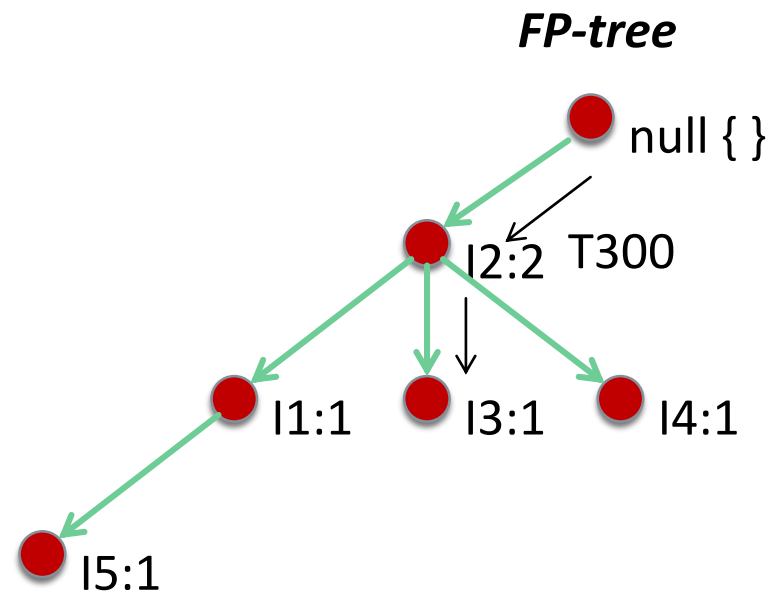


# MINING FREQUENT ITEMSETS

## FP-GROWTH – FP-TREE CONSTRUCTION

$L_1$  - Reordered

Itemset	Support count
{I2}	8
{I1}	7
{I3}	6
{I4}	2
{I5}	2



TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2

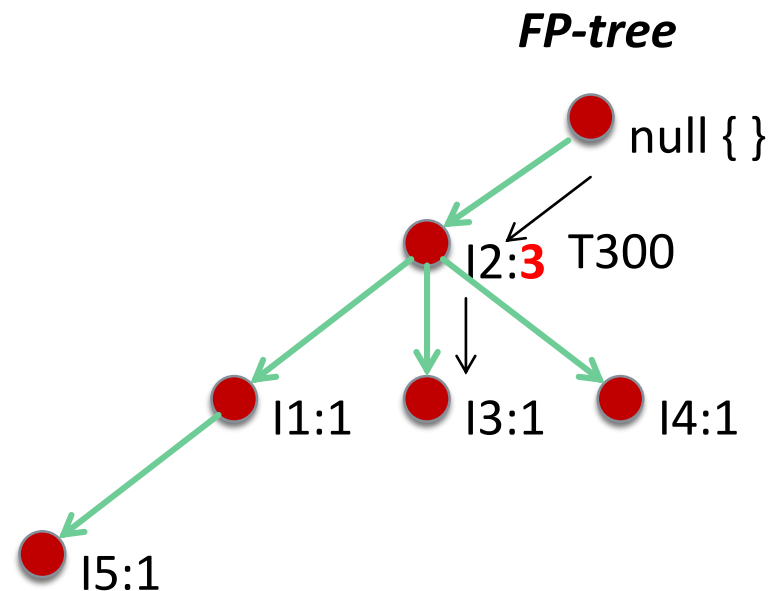


# MINING FREQUENT ITEMSETS

## FP-GROWTH – FP-TREE CONSTRUCTION

$L_1$  - Reordered

Itemset	Support count
{I2}	8
{I1}	7
{I3}	6
{I4}	2
{I5}	2



TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2



# MINING FREQUENT ITEMSETS

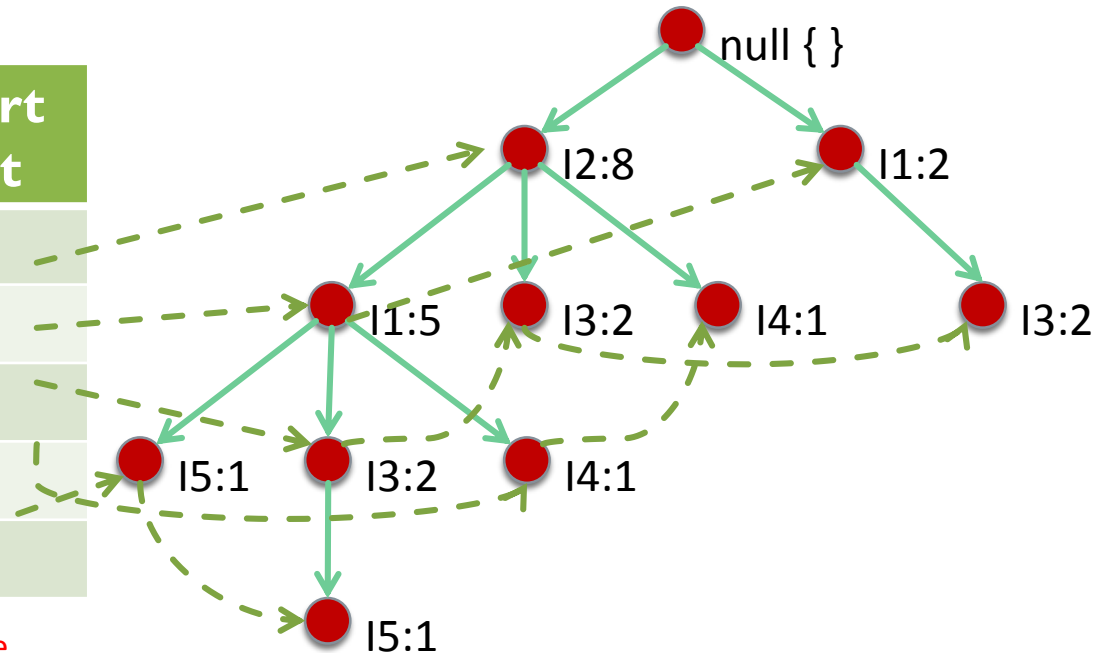
## FP-GROWTH – FP-TREE CONSTRUCTION

$L_1$  - Reordered

Itemset	Support count
{I2}	8
{I1}	7
{I3}	6
{I4}	2
{I5}	2

For Tree Traversal

FP-tree



TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2

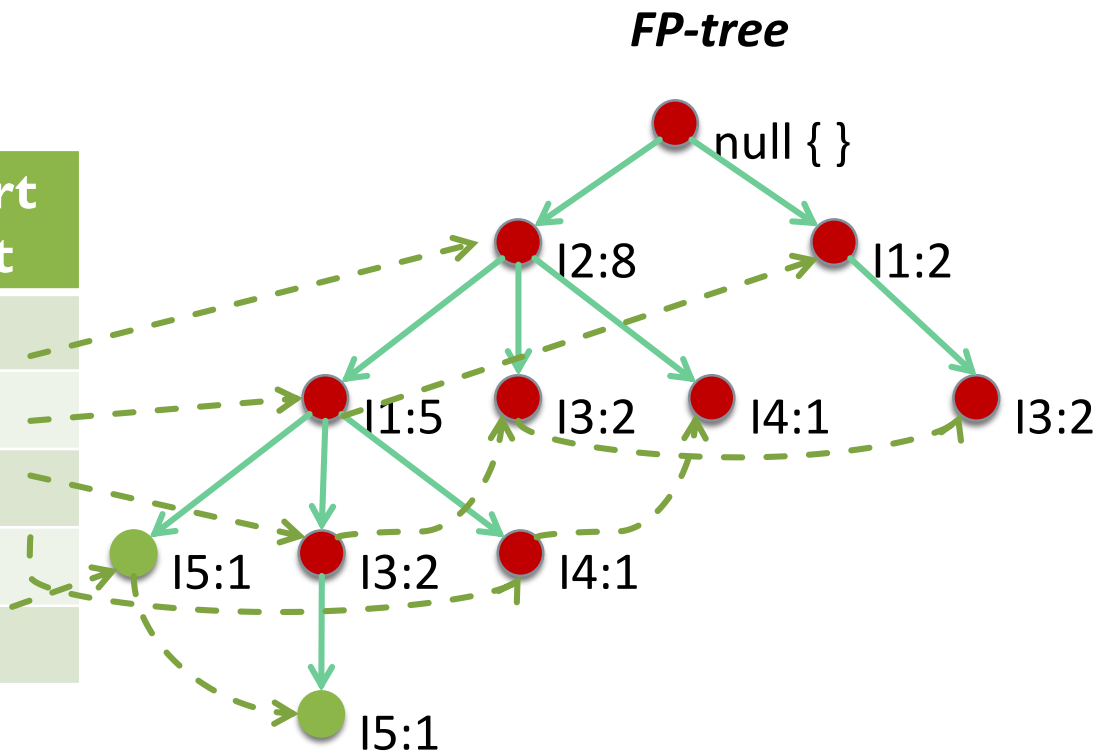


# MINING FREQUENT ITEMSETS

## FP-GROWTH – FREQUENT PATTERNS MINING

$L_1$  - Reordered

Itemset	Support count
{I2}	8
{I1}	7
{I3}	6
{I4}	2
{I5}	2



Bottom-up algorithm – start from leaves and go up to root

TID	List of items
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3
T1000	I1, I2



# MINING FREQUENT ITEMSETS

## FP-GROWTH – CONDITIONAL FP-TREE CONSTRUCTION

For I5

$L_1$  - Reordered

Itemset	Support count	Linked link
{I2}	8	
{I1}	7	
{I3}	6	
{I4}	2	
{I5}	2	

FP-tree

 null { }

Eliminate I5

Eliminate transactions  
not including I5

TID	List of items
T100	I1, I2, <del>I5</del>
<del>T200</del>	<del>I2, I4</del>
<del>T300</del>	<del>I2, I3</del>
<del>T400</del>	<del>I1, I2, I4</del>
<del>T500</del>	<del>I1, I3</del>
<del>T600</del>	<del>I2, I3</del>
<del>T700</del>	<del>I1, I3</del>
T800	I1, I2, I3, <del>I5</del>
<del>T900</del>	<del>I1, I2, I3</del>
<del>T1000</del>	<del>I1, I2</del>



# MINING FREQUENT ITEMSETS

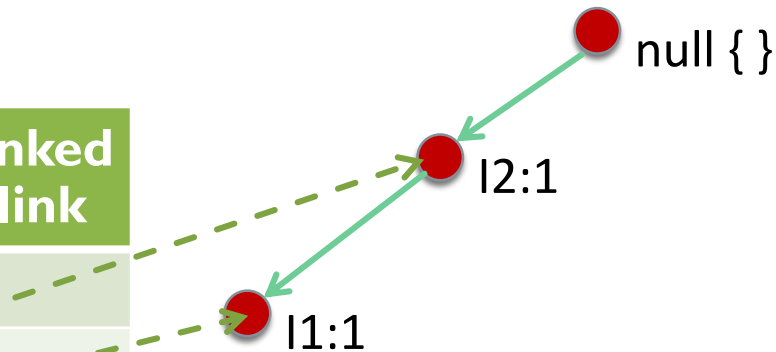
## FP-GROWTH – CONDITIONAL FP-TREE CONSTRUCTION

For I5

$L_1$  - Reordered

Itemset	Support count	Linked link
{I2}	8	
{I1}	7	
{I3}	6	
{I4}	2	
{I5}	2	

FP-tree



Eliminate I5

Eliminate transactions  
not including I5

TID	List of items
T100	I1, I2, <del>I5</del>
<del>T200</del>	<del>I2, I4</del>
<del>T300</del>	<del>I2, I3</del>
<del>T400</del>	<del>I1, I2, I4</del>
<del>T500</del>	<del>I1, I3</del>
<del>T600</del>	<del>I2, I3</del>
<del>T700</del>	<del>I1, I3</del>
T800	I1, I2, I3, <del>I5</del>
<del>T900</del>	<del>I1, I2, I3</del>
<del>T1000</del>	<del>I1, I2</del>



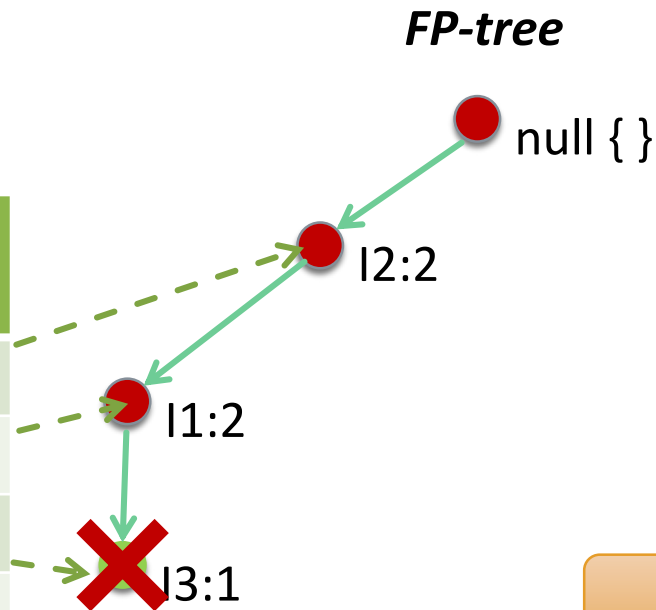
# MINING FREQUENT ITEMSETS

## FP-GROWTH – CONDITIONAL FP-TREE CONSTRUCTION

For I5

$L_1$  - Reordered

Itemset	Support count	Linked link
{I2}	8	
{I1}	7	
{I3}	6	
{I4}	2	
{I5}	2	



Eliminate I5

Eliminate transactions  
not including I5

TID	List of items
T100	I1, I2, <del>I5</del>
<del>T200</del>	<del>I2, I4</del>
<del>T300</del>	<del>I2, I3</del>
<del>T400</del>	<del>I1, I2, I4</del>
<del>T500</del>	<del>I1, I3</del>
<del>T600</del>	<del>I2, I3</del>
<del>T700</del>	<del>I1, I3</del>
T800	I1, I2, I3, <del>I5</del>
<del>T900</del>	<del>I1, I2, I3</del>
<del>T1000</del>	<del>I1, I2</del>





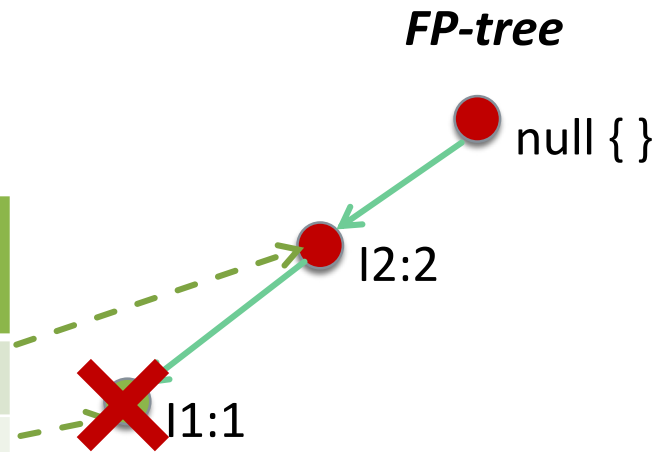
# MINING FREQUENT ITEMSETS

## FP-GROWTH – CONDITIONAL FP-TREE CONSTRUCTION

For I4

*L<sub>1</sub> - Reordered*

Itemset	Support count	Linked link
{I2}	8	
{I1}	7	
{I3}	6	
{I4}	2	
{I5}	2	



Eliminate I4

Eliminate transactions  
not including I4

TID	List of items
<del>T100</del>	<del>I1, I2, I5</del>
T200	I2, <del>I4</del>
<del>T300</del>	<del>I2, I3</del>
T400	I1, I2, <del>I4</del>
<del>T500</del>	<del>I1, I3</del>
<del>T600</del>	<del>I2, I3</del>
<del>T700</del>	<del>I1, I3</del>
<del>T800</del>	<del>I1, I2, I3, I5</del>
<del>T900</del>	<del>I1, I2, I3</del>
<del>T1000</del>	<del>I1, I2</del>



# MINING FREQUENT ITEMSETS

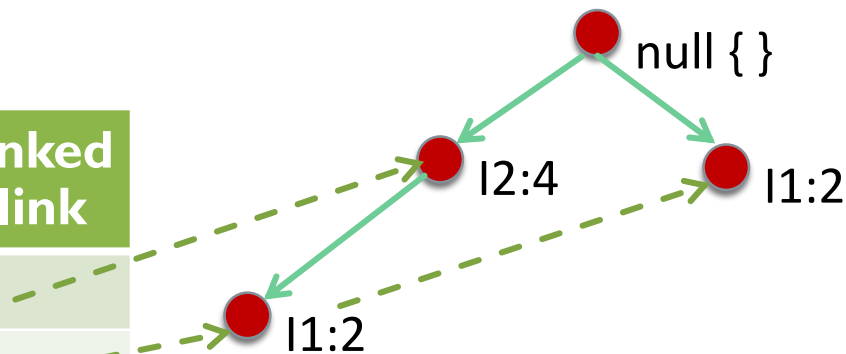
## FP-GROWTH – CONDITIONAL FP-TREE CONSTRUCTION

For I3

$L_1$  - Reordered

Itemset	Support count	Linked link
{I2}	8	
{I1}	7	
{I3}	6	
{I4}	2	
{I5}	2	

FP-tree



Eliminate I3

Eliminate transactions  
not including I3

TID	List of items
<del>T100</del>	<del>I1, I2, I5</del>
<del>T200</del>	<del>I2, I4</del>
T300	I2, <del>I3</del>
<del>T400</del>	<del>I1, I2, I4</del>
T500	I1, <del>I3</del>
T600	I2, <del>I3</del>
T700	I1, <del>I3</del>
T800	I1, I2, <del>I3</del> , I5
T900	I1, I2, <del>I3</del>
<del>T1000</del>	<del>I1, I2</del>



# MINING FREQUENT ITEMSETS

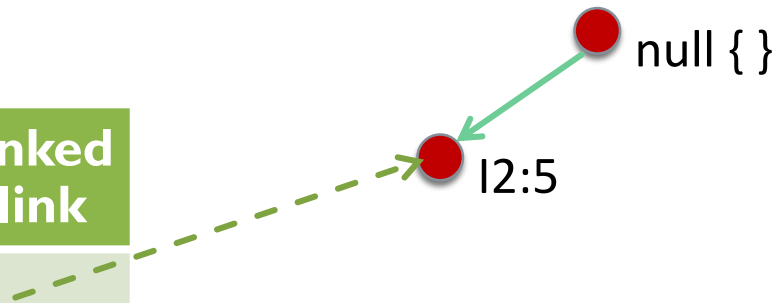
## FP-GROWTH – CONDITIONAL FP-TREE CONSTRUCTION

For I1

*L<sub>1</sub> - Reordered*

Itemset	Support count	Linked link
{I2}	8	
{I1}	7	
{I3}	6	
{I4}	2	
{I5}	2	

*FP-tree*



Eliminate I1

Eliminate transactions  
not including I1

TID	List of items
T100	<del>I1</del> , I2, I5
<del>T200</del>	<del>I2, I4</del>
<del>T300</del>	<del>I2, I3</del>
T400	<del>I1</del> , I2, I4
T500	<del>I1</del> , I3
<del>T600</del>	<del>I2, I3</del>
T700	<del>I1</del> , I3
T800	<del>I1</del> , I2, I3, I5
T900	<del>I1</del> , I2, I3
T1000	<del>I1</del> , I2

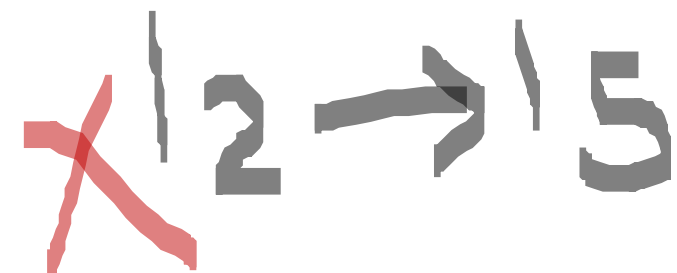
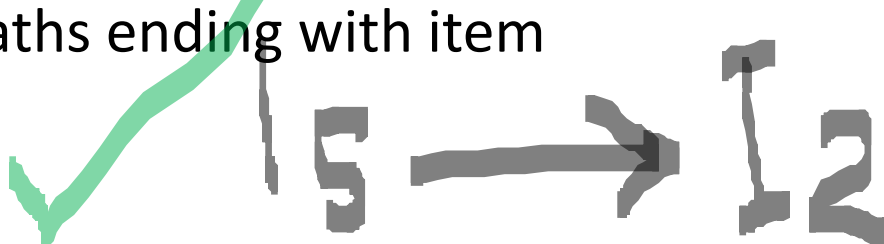


# MINING FREQUENT ITEMSETS

## FP-GROWTH

Item	Conditional Pattern Base	Conditional FP-tree	Frequent Patterns Generated
15	$\{\{I2, I1: 1\}, \{I2, I1, I3: 1\}\}$	$\langle I2:2, I1:2 \rangle$	$\{I2, I5: 2\}, \{I1, I5: 2\}, \{I2, I1, I5: 2\}$
14	$\{\{I2, I1: 1\}, \{I2: 1\}\}$	$\langle I2:2 \rangle$	$\{I2, I4: 2\}$
13	$\{\{I2, I1: 2\}, \{I2: 2\}, \{I1: 2\}\}$	$\langle I2:4, I1:2 \rangle, \langle I1:2 \rangle$	$\{I2, I3: 4\}, \{I1, I3: 4\}, \{I2, I1, I3: 2\}$
11	$\{\{I2: 5\}\}$	$\langle I2:5 \rangle$	$\{I2, I1: 5\}$

Paths ending with item



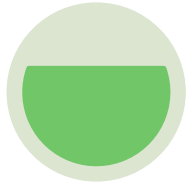


## The Basics

Market Basket  
Analysis

Frequent Itemsets

Association Rules

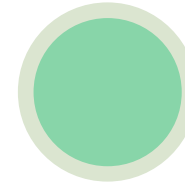


## Frequent Itemset Mining Methods

Apriori Algorithm

Generating  
Association Rules  
from Frequent  
Itemsets

FP-Growth



## Pattern Evaluation Methods



# PATTERN EVALUATION METHODS

- Not all association rules are interesting
  - $\text{Buys}(X, \text{"Computer games"}) \rightarrow \text{buys}(X, \text{"Videos"})$  [40%, 66%]
  - $P(\text{"videos"})$  is 75% > 66%
  - The two items are negatively associated means buying one decreases the likelihood of buying the other
  - We need to measure “real strength” of rule
- **Correlation analysis**
  - $A \rightarrow B$  [support , confidence , correlation]



# PATTERN EVALUATION METHODS

1. **Lift** =  $\frac{P(A \cup B)}{P(A)P(B)}$ 
  - A and B are independent if  $P(A \cup B) = P(A)P(B)$
  - Otherwise, **dependent and correlated** occurrence
  - If lift < 1, A is **Negatively correlated** with B
  - If lift > 1, A is **Positively correlated** with B ..... A's occurrence "lifts" the occurrence of B
2.  $\chi^2$  → already discussed in previous lecture



---

# QUESTIONS?

NEXT ...

