

关联规则挖掘

丁兆云

关联规则挖掘

丁兆云

主要内容

- ◆ 1. 基于概念
- ◆ 2. 频繁项挖掘算法
- ◆ 3. 关联分析的评估

1.1 定义: 关联分析 (association analysis)

- ◆ 关联分析用于发现隐藏在大型数据集中的令人感兴趣的联系，所发现的模式通常用**关联规则**或**频繁项集**的形式表示。
- ◆ 关联规则反映一个**事物与其它事物**之间的相互依存性和关联性。如果**两个或者多个事物之间**存在一定的关联关系，那么，其中一个事物发生就能够预测与它相关联的其它事物的发生。

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Rules Discovered:
 $\{\text{Diaper}\} \rightarrow \{\text{Beer}\}$

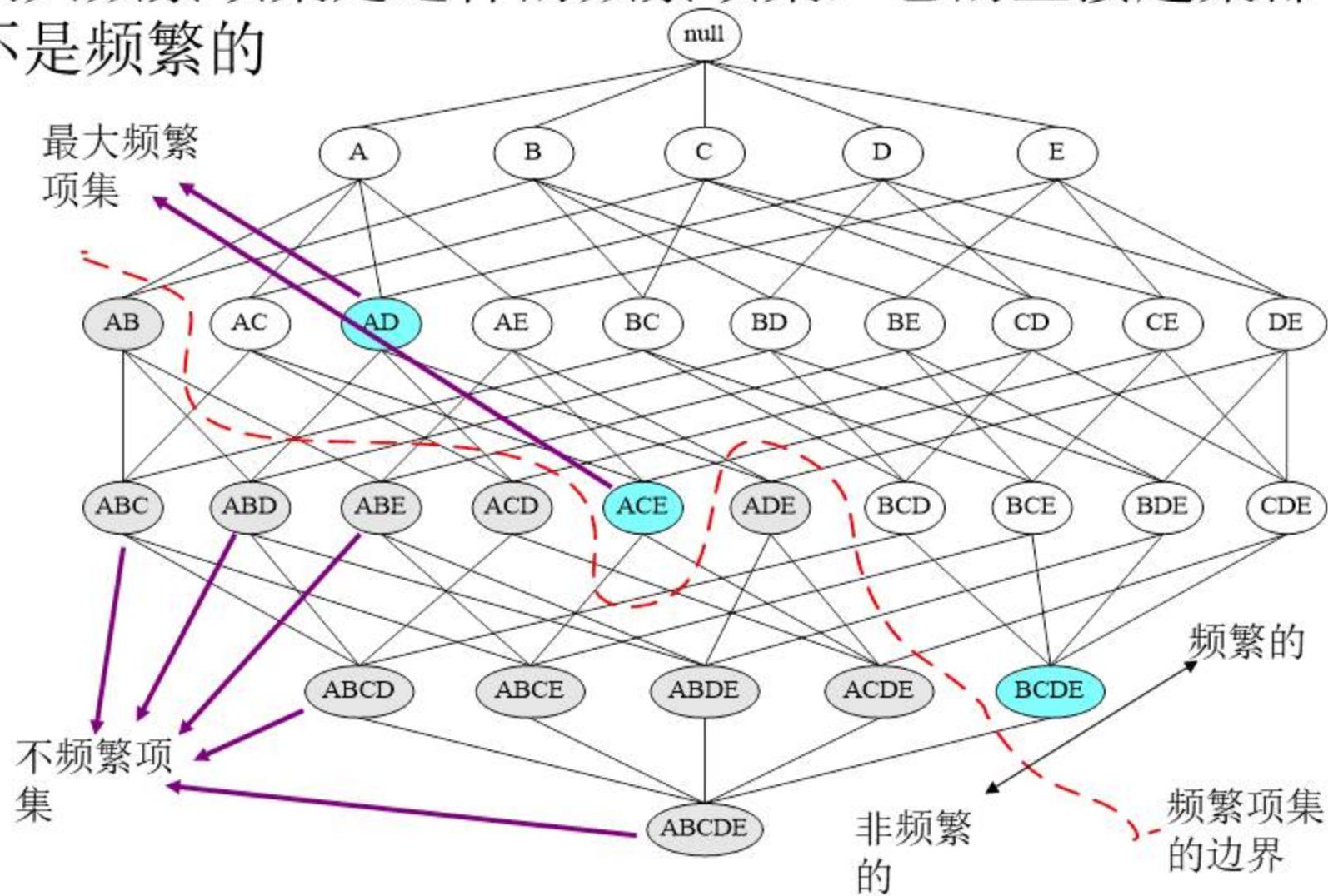
1.1 定义: 频繁项集 (Frequent Itemset)

- 项集 (Itemset)
 - 包含0个或多个项的集合
例子: {Milk, Bread, Diaper}
 - k-项集
如果一个项集包含k个项
- 支持度计数 (Support count) (σ)
 - 包含特定项集的事务个数
 - 例如: $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$
- 支持度 (Support)
 - 包含项集的事务数与总事务数的比值
 - 例如: $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$
- 频繁项集 (Frequent Itemset)
 - 满足最小支持度阈值 ($minsup$) 的所有项集

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

1.2最大频繁项集（Maximal Frequent Itemset）

最大频繁项集是这样的频繁项集，它的直接超集都不是频繁的



1.3 定义: 关联规则 (Association Rule)

- 关联规则

- 关联规则是形如 $X \rightarrow Y$ 的蕴含表达式, 其中 X 和 Y 是不相交的项集
- 例子:
 $\{Milk, Diaper\} \rightarrow \{Beer\}$

- 关联规则的强度

- 支持度 Support (s)
 - ◆ 确定项集的频繁程度
- 置信度 Confidence (c)
 - ◆ 确定 Y 在包含 X 的事务中出现的频繁程度

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example:

$$\{Milk, Diaper\} \Rightarrow Beer$$

$$s = \frac{\sigma(Milk, Diaper, Beer)}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(Milk, Diaper, Beer)}{\sigma(Milk, Diaper)} = \frac{2}{3} = 0.67$$

1.4 关联规则挖掘问题

- ◆ 关联规则挖掘问题：给定事务的集合 T , 关联规则发现是指找出支持度大于等于 $minsup$ 并且置信度大于等于 $minconf$ 的所有规则, $minsup$ 和 $minconf$ 是对应的支持度和置信度阈值

事务Id	项集合
10	A, B, C
20	A, C
30	A, D
40	B, E, F

最小支持度 50%
最小置信度 50%

频繁模式	支持度
{A}	75%
{B}	50%
{C}	50%
{A, C}	50%

$$\text{支持度} = \text{support}(\{A\} \cup \{C\}) = \text{[填空1]}$$

$$\text{置信度} = \text{support}(\{A\} \cup \{C\}) / \text{support}(\{A\}) = \text{[填空2]}$$

正常使用填空题需3.0以上版本雨课堂

作答

1.4 关联规则挖掘问题

- 关联规则挖掘问题：给定事务的集合 T , 关联规则发现是指找出支持度大于等于 $minsup$ 并且置信度大于等于 $minconf$ 的所有规则, $minsup$ 和 $minconf$ 是对应的支持度和置信度阈值

事务Id	项集合
10	A, B, C
20	A, C
30	A, D
40	B, E, F

最小支持度 50%
最小置信度 50%

频繁模式	支持度
{A}	75%
{B}	50%
{C}	50%
{A, C}	50%

规则 $A \Rightarrow C$:

$$\text{支持度} = \text{support}(\{A\} \cup \{C\}) = 50\%$$

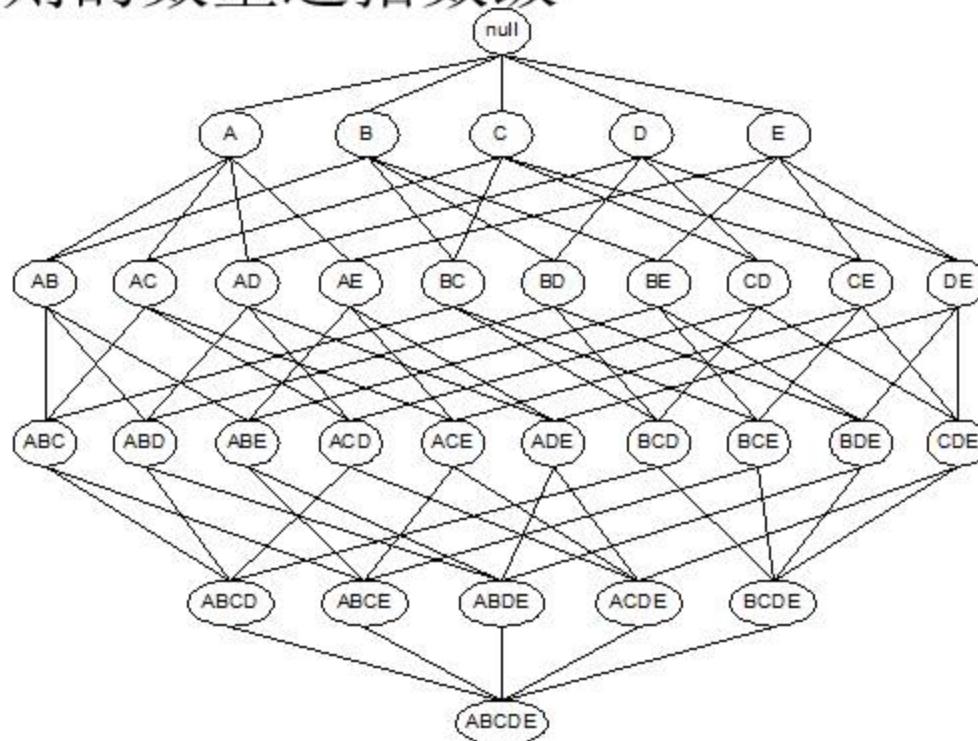
$$\text{置信度} = \text{support}(\{A\} \cup \{C\}) / \text{support}(\{A\}) = 66.6\%$$

1.4 挖掘关联规则 (Mining Association Rules)

- ◆ 大多数关联规则挖掘算法通常采用的一种策略是，将关联规则挖掘任务分解为如下两个主要的子任务：
 1. 频繁项集产生 (Frequent Itemset Generation)
 - 其目标是发现满足最小支持度阈值的所有项集，这些项集称作频繁项集。
 2. 规则的产生 (Rule Generation)
 - 其目标是从上一步发现的频繁项集中提取所有高置信度的规则，这些规则称作强规则 (strong rule)。

1.5 关联规则原始方法

- ◆ 挖掘关联规则的一种原始方法是：Brute-force approach：
 - 计算每个可能规则的支持度和置信度
 - 这种方法计算代价过高，因为可以从数据集提取的规则的数量达指数级



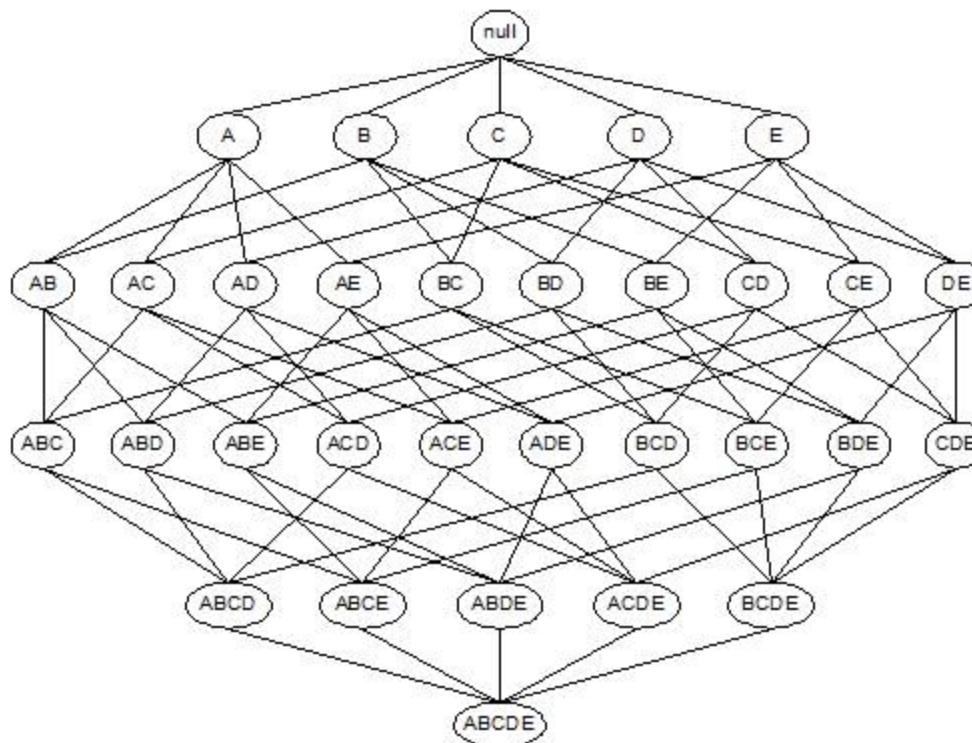
主要内容

- ◆ 1. 基于概念
- ◆ 2. 频繁项挖掘算法
- ◆ 3. 关联分析的评估

2频繁项集产生 (Frequent Itemset Generation)

- ◆ Brute-force 方法:

- 计算开销大



2降低产生频繁项集计算复杂度的方法

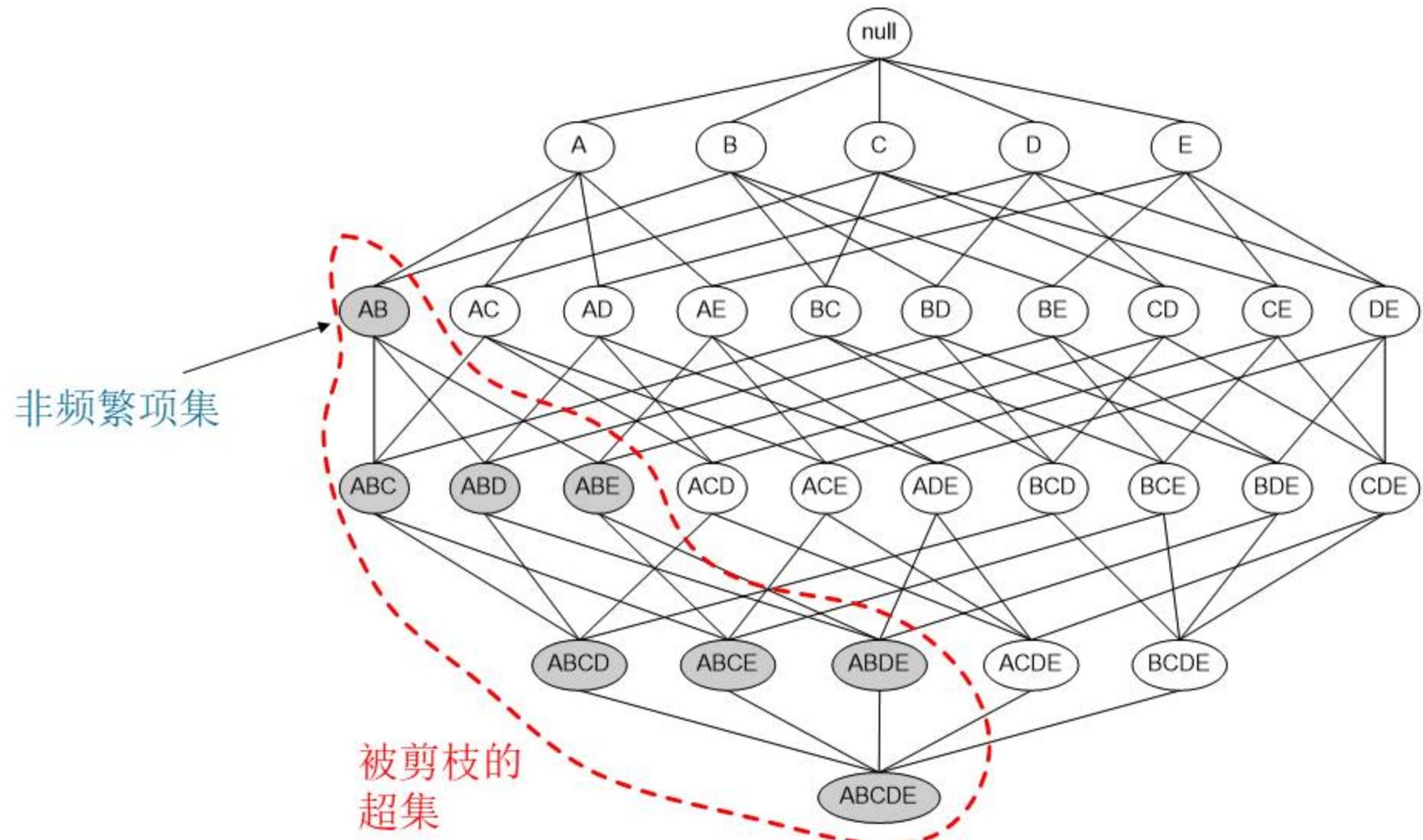
- ◆ 减少候选项集的数量
 - 先验原理:(Apriori)
- ◆ 减少比较的次数
 - 替代将每个候选项集与每个事务相匹配，可以使用更高级的数据结构，或存储候选项集或压缩数据集，来减少比较次数(FPGrowth)

2.1 Apriori

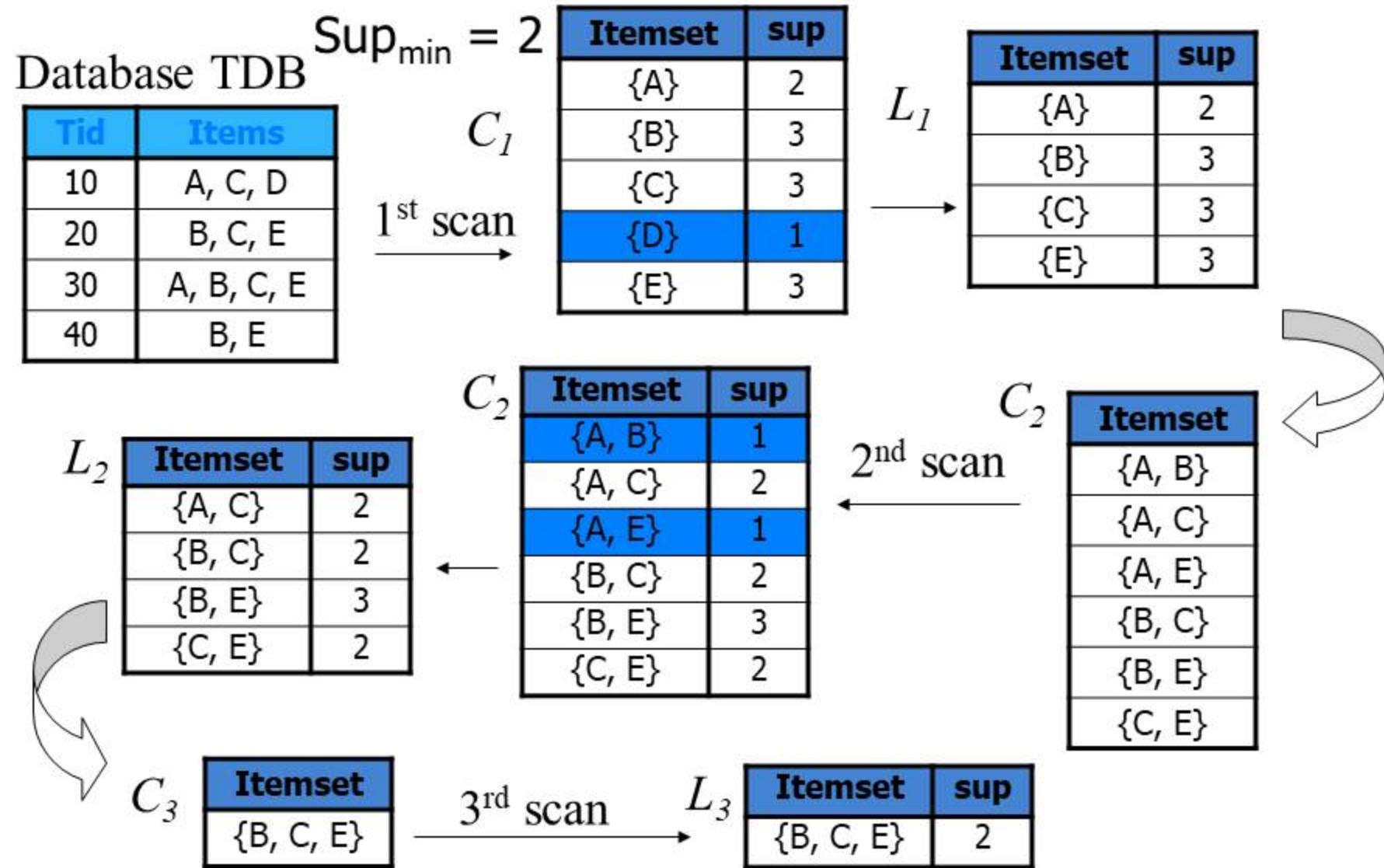
- ◆ 先验原理：

- 如果一个项集是频繁的，则它的所有子集一定也是频繁的
- 相反，如果一个项集是非频繁的，则它的所有超集也一定是非频繁的：

例子



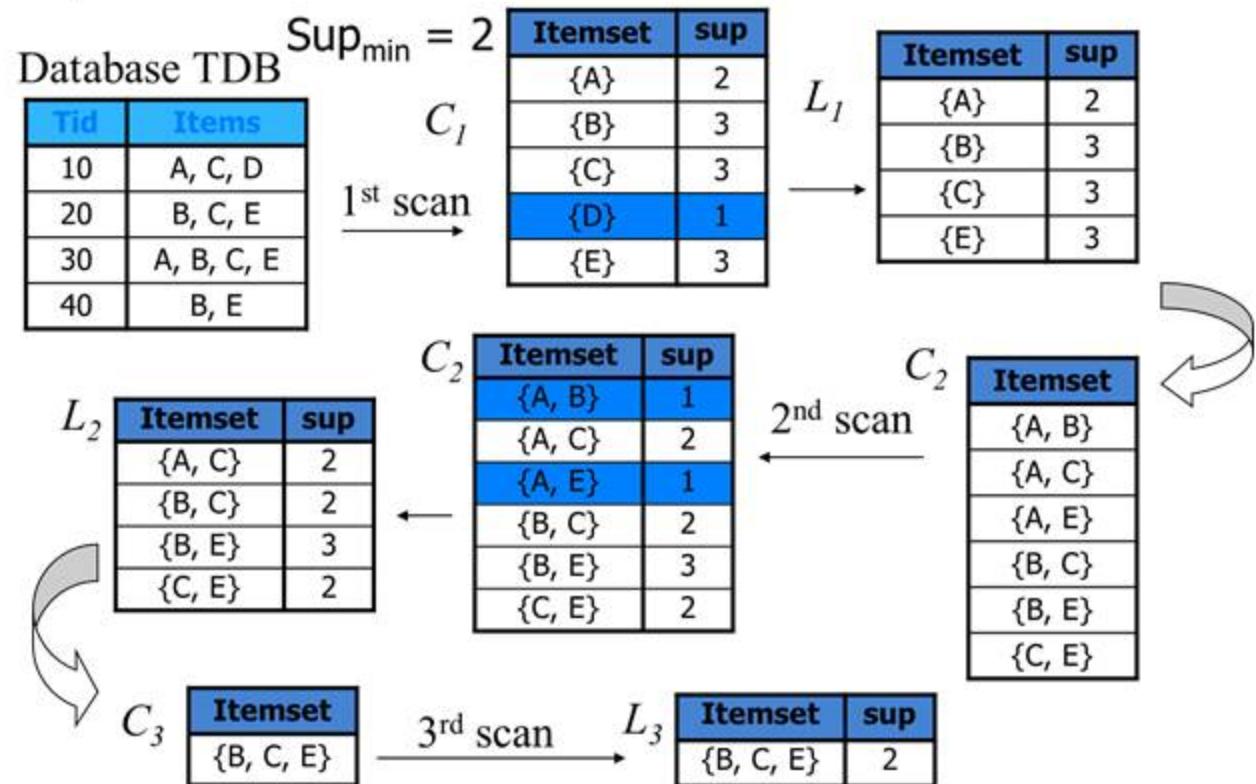
2.1 Apriori算法过程：最小支持度计数=2



2.1 Apriori算法注意问题-项的字典序

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

尽管集合具有无序性，但为了方便比较计数，通常对所有商品做一个默认的排序（类似于建立一个字典索引）



2.1 Apriori算法注意问题-项的连接

对于任何2个需要连接的项集 A, B, C
B, C, E

2.1 Apriori算法注意问题-项的连接

对于任何2个需要连接的项集 A, B, C 去掉第1个项集的首项
B, C, E

2.1 Apriori算法注意问题-项的连接

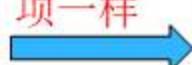
对于任何2个需要连接的项集 A, B, C 去掉第1个项集的首项
 B, C, E 去掉第2个项集的尾项

2.1 Apriori算法注意问题-项的连接

对于任何2个需要连接的项集

A, B, C 去掉第1个项集的首项
B, C, E 去掉第2个项集的尾项

若剩下的
项一样



B, C

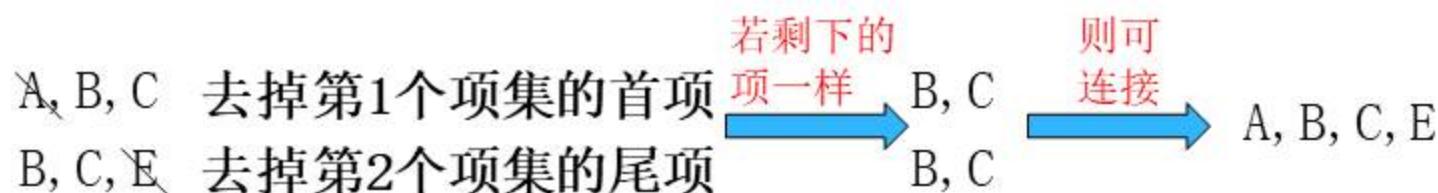
2.1 Apriori算法注意问题-项的连接

对于任何2个需要连接的项集

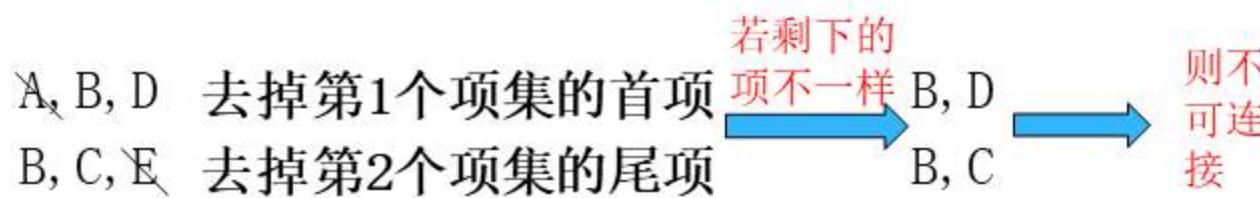


2.1 Apriori算法注意问题-项的连接

对于任何2个需要连接的项集



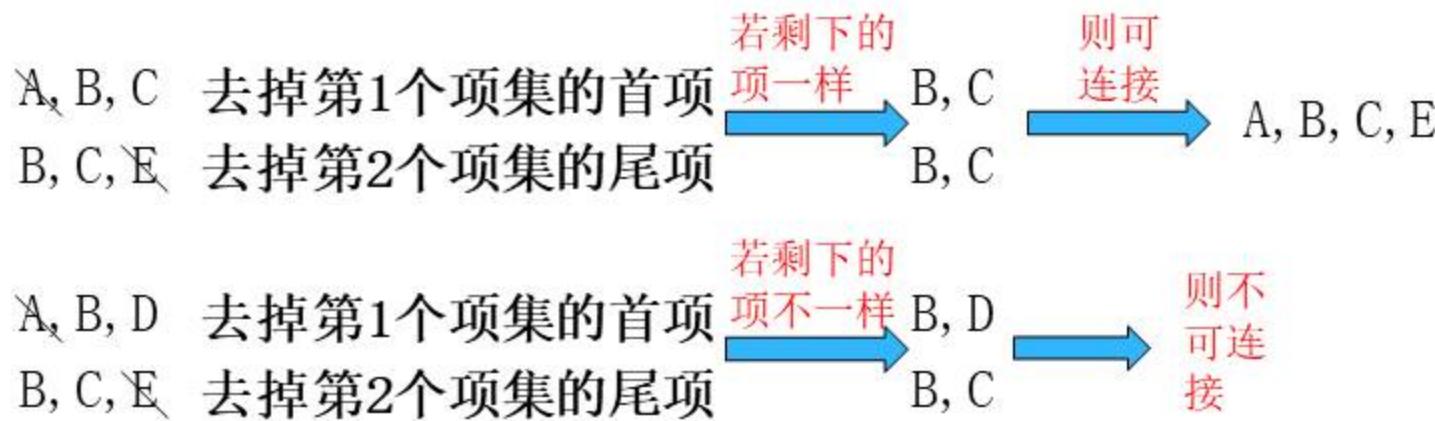
对于任何2个需要连接的项集



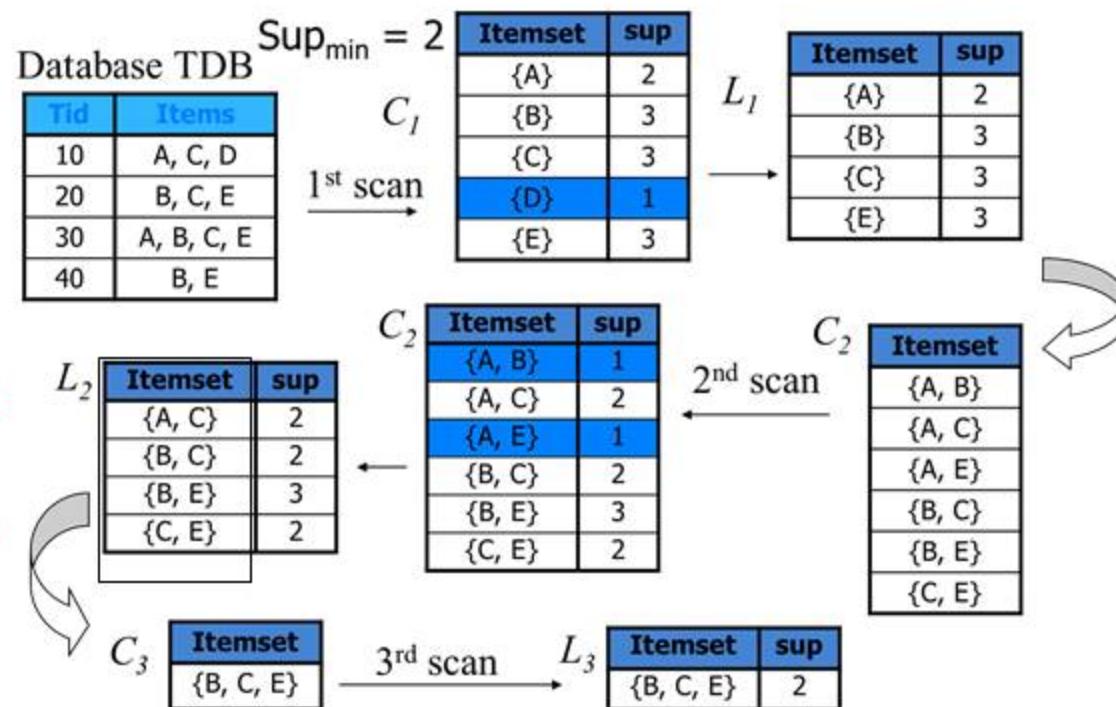
2.1 Apriori算法注意问题-项的连接

对于任何2个需要连接的项集

对于任何2个需要连接的项集

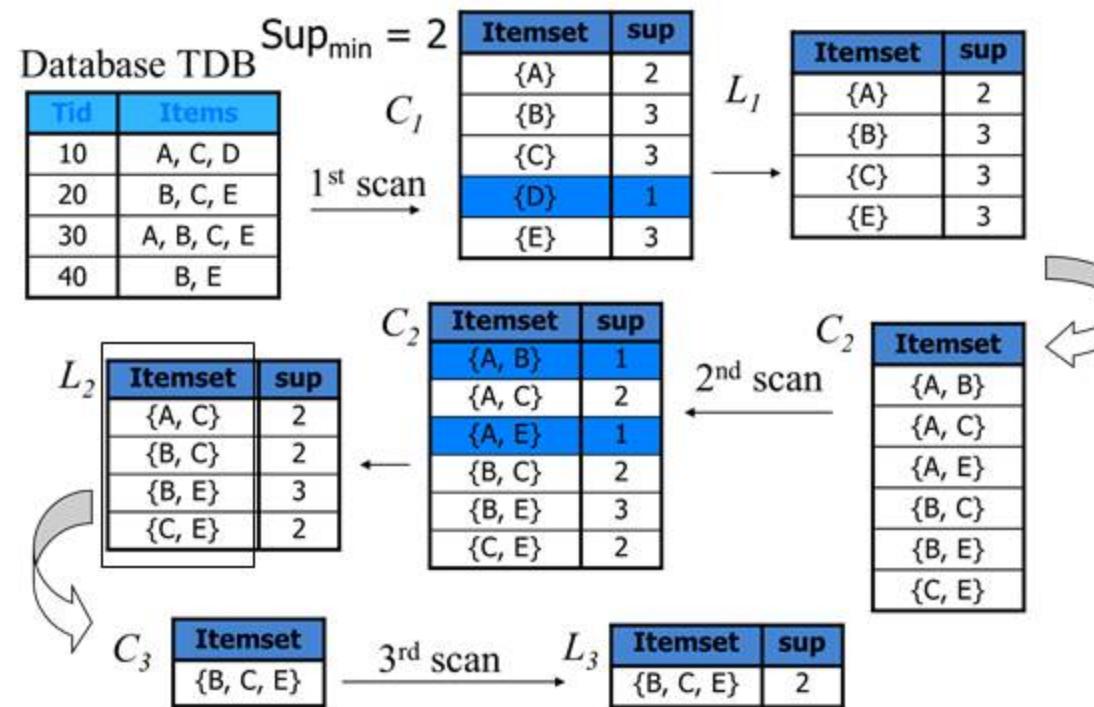


频繁2项集生成的时候选3项集是什么？



频繁2项集生成的候选3项集是什么？

- A A,B,C
- B A,C,E
- C A,B,E
- D B,C,E

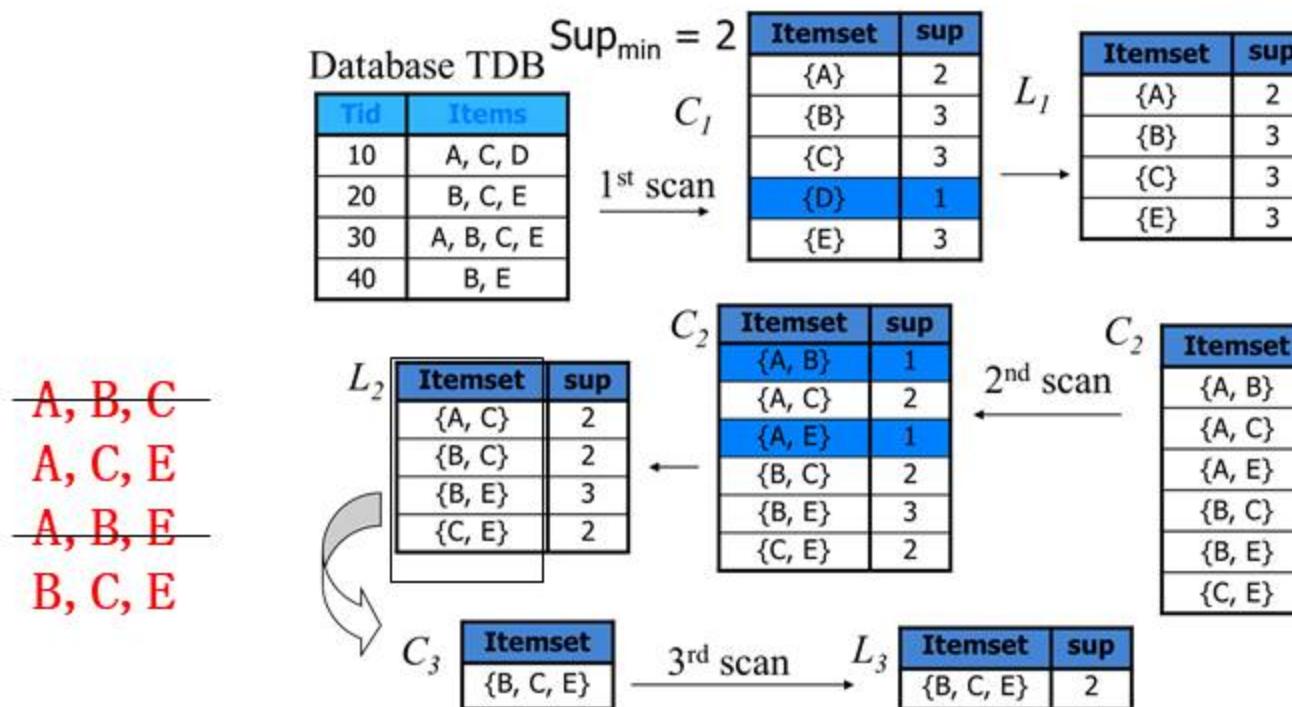
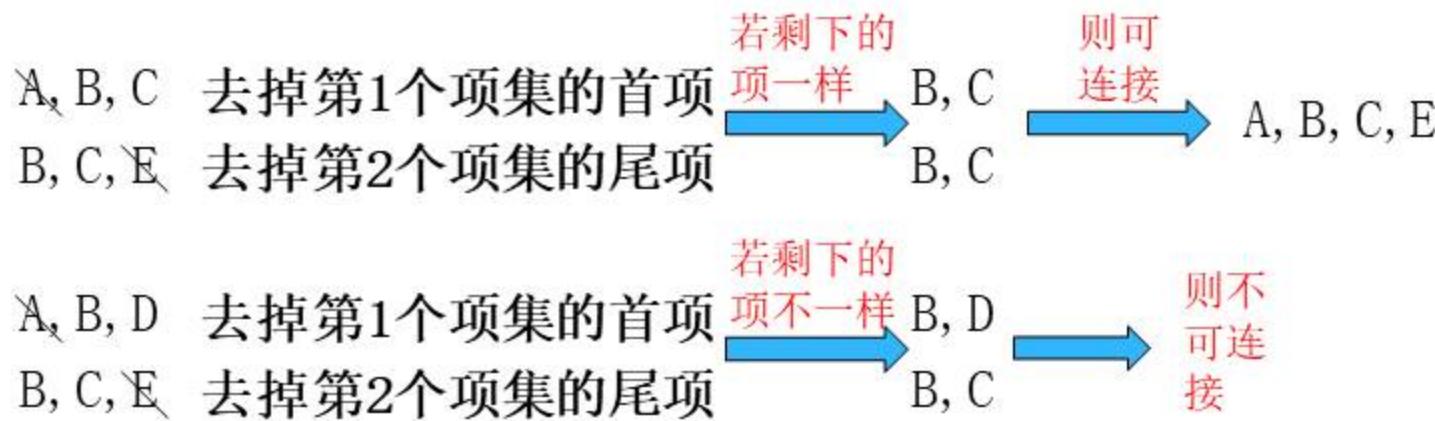


提交

2.1 Apriori算法注意问题-项的连接

对于任何2个需要连接的项集

对于任何2个需要连接的项集

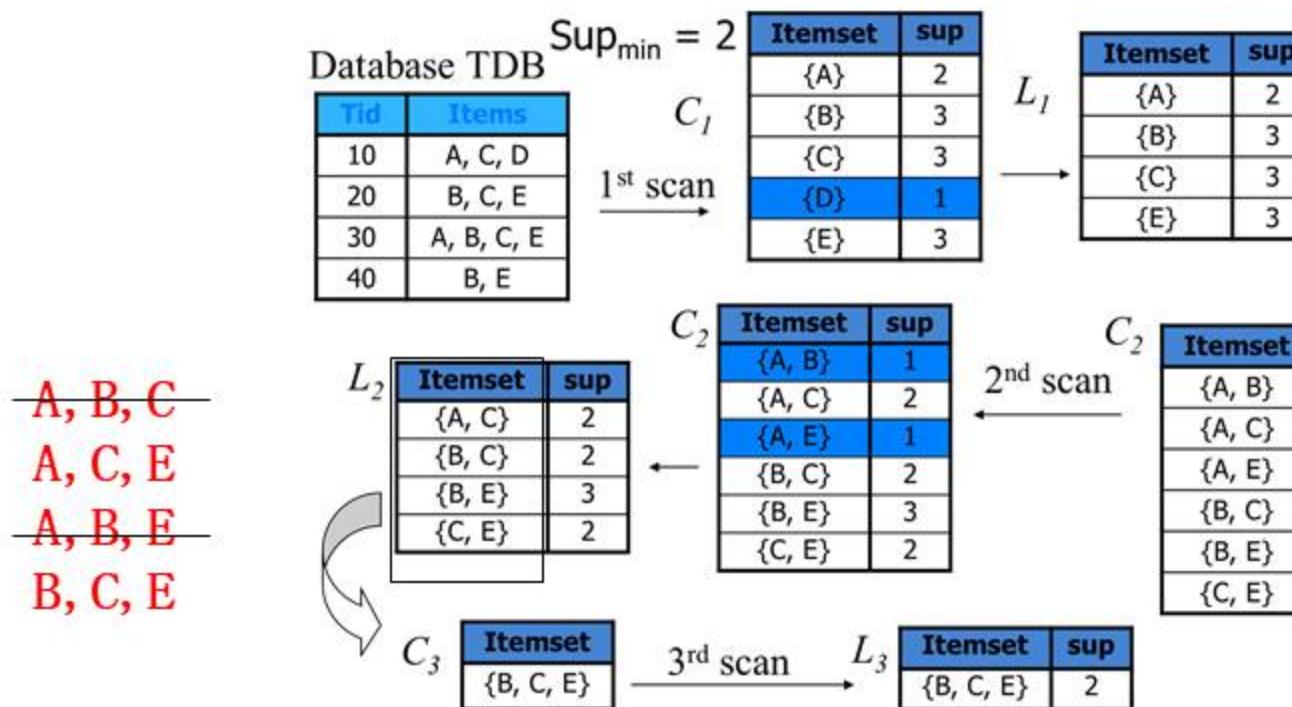
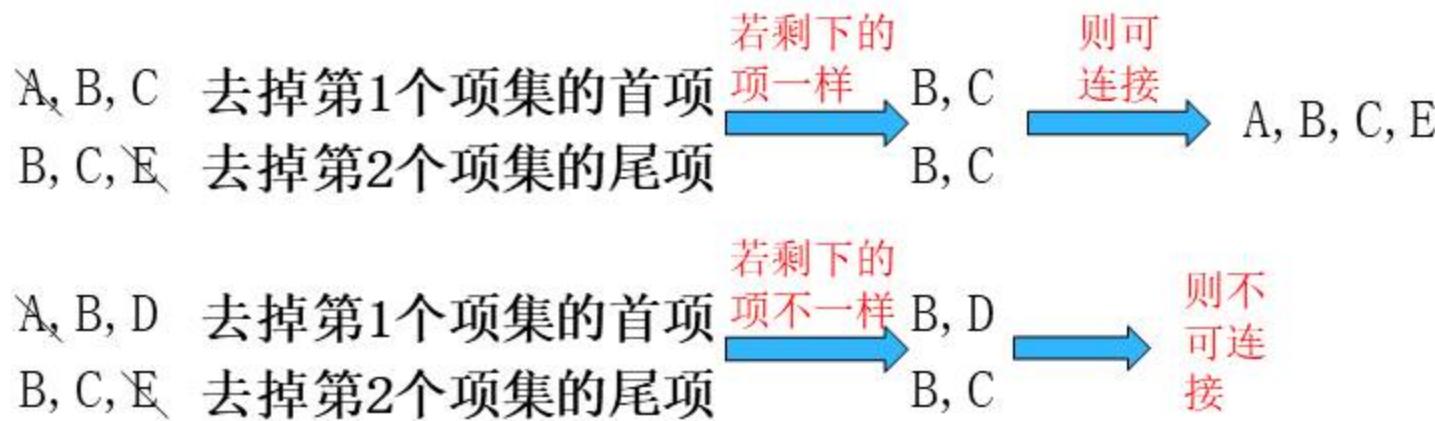


2.1 Apriori算法注意问题-项的连接

对于任何2个需要连接的项集

对于任何2个需要连接的项集

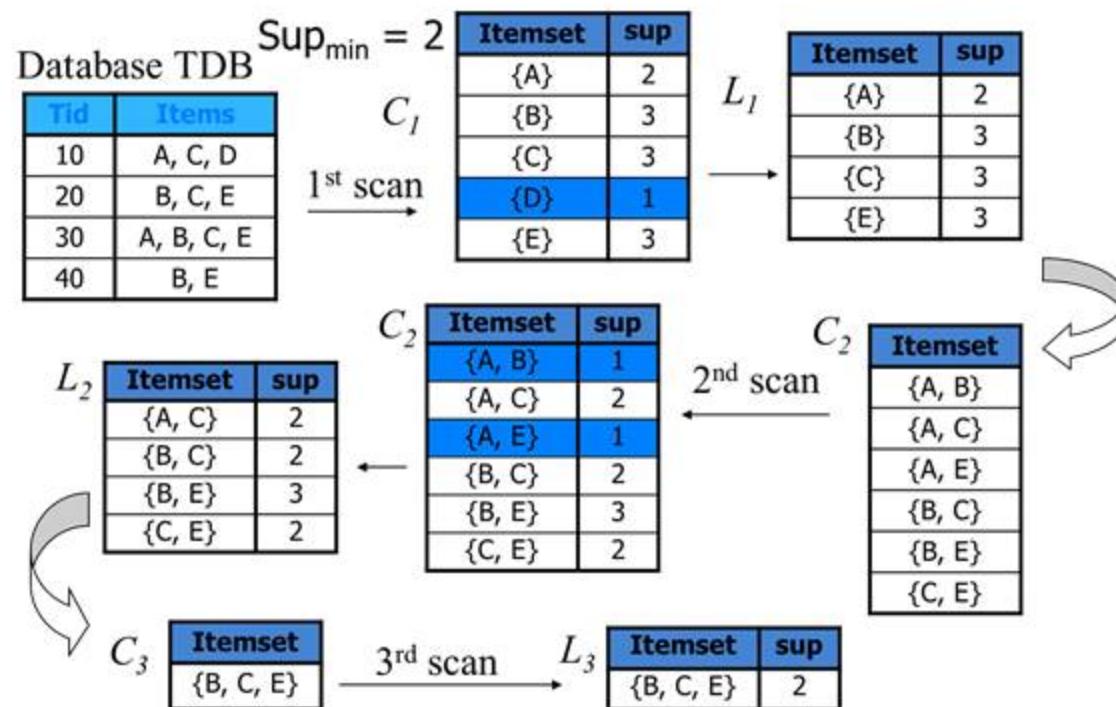
项的连接准则的优点是降低候选项的生成



A, B, C
A, C, E
A, B, E
B, C, E

2.1 Apriori 算法特点

- ◆ 多次扫描数据库
- ◆ 候选项规模庞大
- ◆ 计算支持度开销大

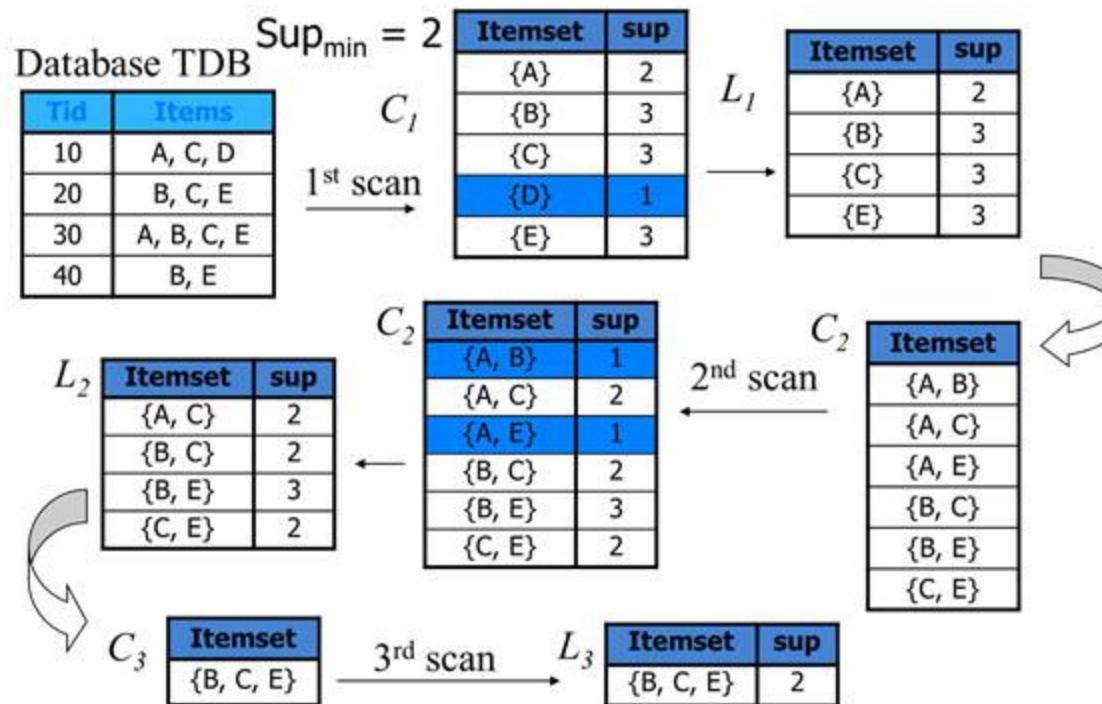


2.1 提高Apriori算法性能的方法

- ◆ Hash-based itemset counting (散列项集计数)
- ◆ Transaction reduction (事务压缩)
- ◆ Partitioning (划分)
- ◆ Sampling (采样)

2.1 Apriori算法缺点

- Apriori算法需要反复的生成候选项，如果项的数目比较大，候选项的数目将达到**组合爆炸式**的增长



2.2 FP Growth

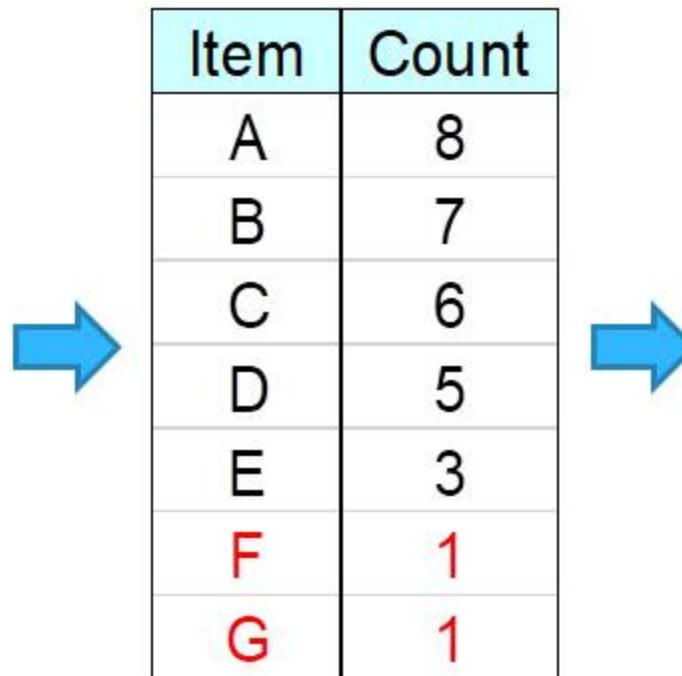
- ◆ 背景
 - 韩家炜等人，2000年
- ◆ 基本思想
 - 只扫描数据库两遍，构造频繁模式树（FP-Tree）
 - 自底向上递归产生频繁项集
 - FP树是一种输入数据的压缩表示，它通过逐个读入事务，并**把每个事务映射到FP树中的一条路径来构造。**

2.2构造FP树：第一遍扫描

原始事物

TID	Items
1	{B,F,A}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{A}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}
11	{G}

删除支持度小于2的项



The diagram illustrates the process of constructing an FP tree. It shows three tables connected by blue arrows. The first arrow points from the '原始事物' (Original Transactions) table to the second table. The second arrow points from the second table to the third table.

Item	Count
A	8
B	7
C	6
D	5
E	3
F	1
G	1

Items按照出现次数降序排列

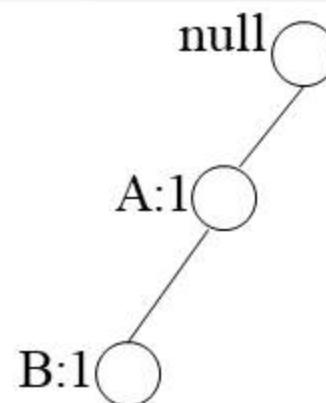
TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{A}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

假设事先指定最小支持度计数为2

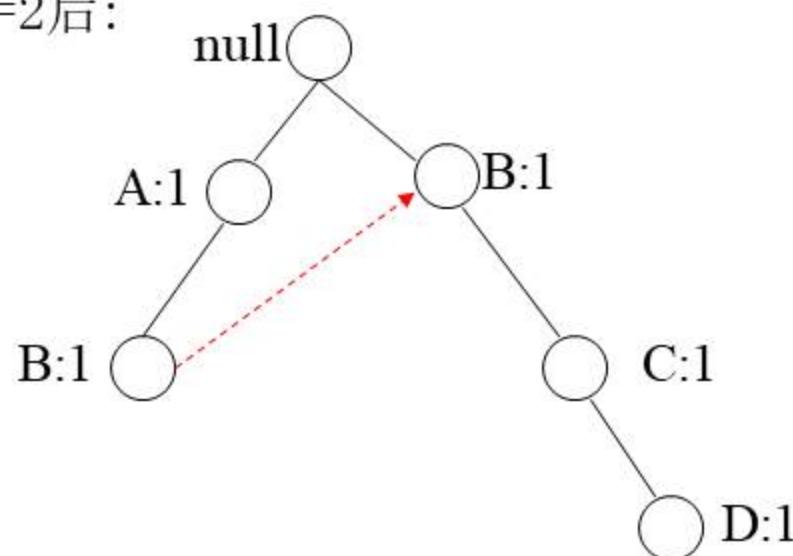
2.2构造FP树：第二遍扫描

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{A}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

读入事务 TID=1后：



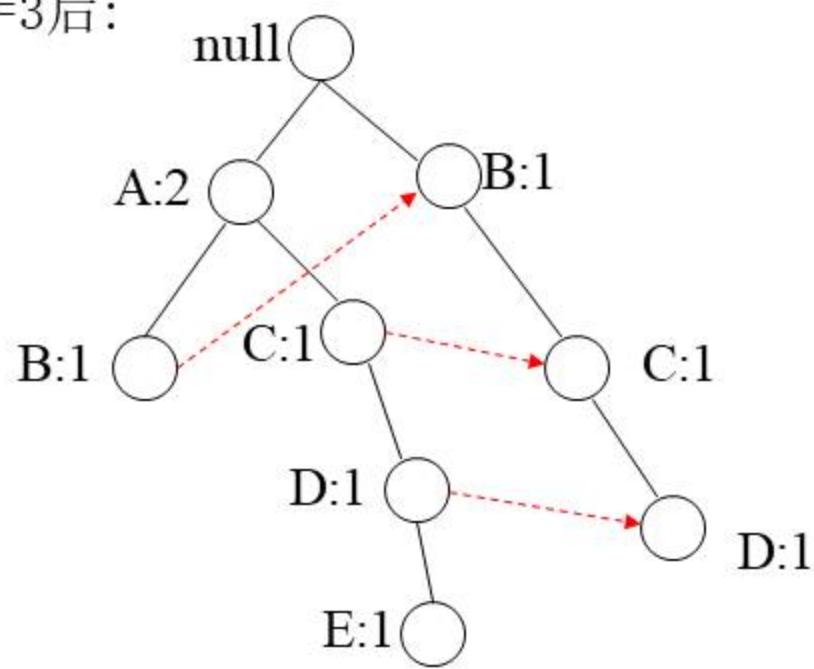
读入事务 TID=2后：



2.2构造FP树：第二遍扫描

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{A}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

读入事务 TID=3后：

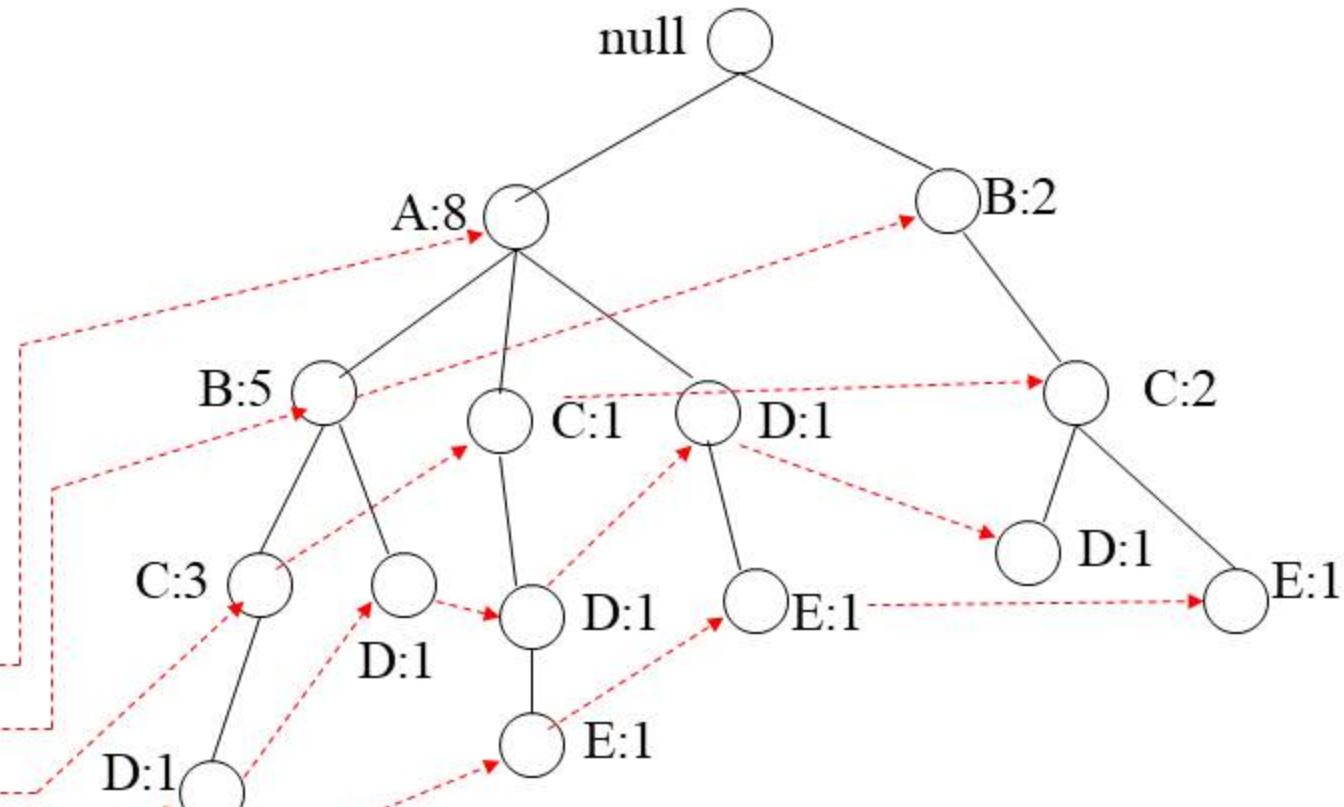


2.2构造FP树：第二遍扫描

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{A}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Header table

Item	Pointer
A	-
B	-
C	-
D	-
E	-



2.2构造FP树：用FP-tree挖掘频繁集

- ◆ 基本思想(分治)
 - 用FP-tree递归增长频繁集
- ◆ 方法
 - 对每个项，生成它的**条件模式基**，然后生成它的**条件 FP-tree**
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为**空**，或只含**唯一的一个路径**(此路径的每个子路径对应的项集都是频繁集)

练习：构造FP树

<u>TID</u>	<u>Items bought</u>	<u>(ordered) frequent items</u>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

最小支持度 = 0.5

步骤：

1. 扫描数据库一次，得到频繁1-项集
2. 把项按支持度递减排序
3. 再一次扫描数据库，建立FP-tree

头表

<u>Item frequency head</u>
f 4
c 4
a 3
b 3
m 3
p 3

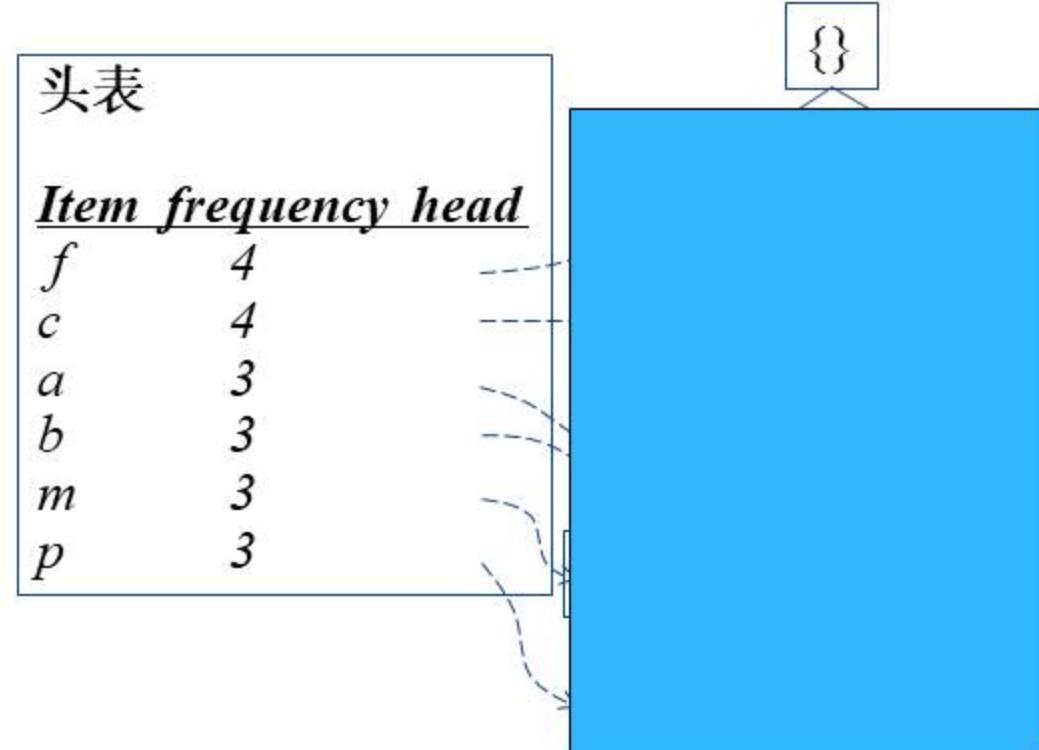
练习：构造FP树

<u>TID</u>	<u>Items bought</u>	<u>(ordered) frequent items</u>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

最小支持度 = 0.5

步骤：

1. 扫描数据库一次，得到频繁1-项集
2. 把项按支持度递减排序
3. 再一次扫描数据库，建立FP-tree



练习：构造FP树

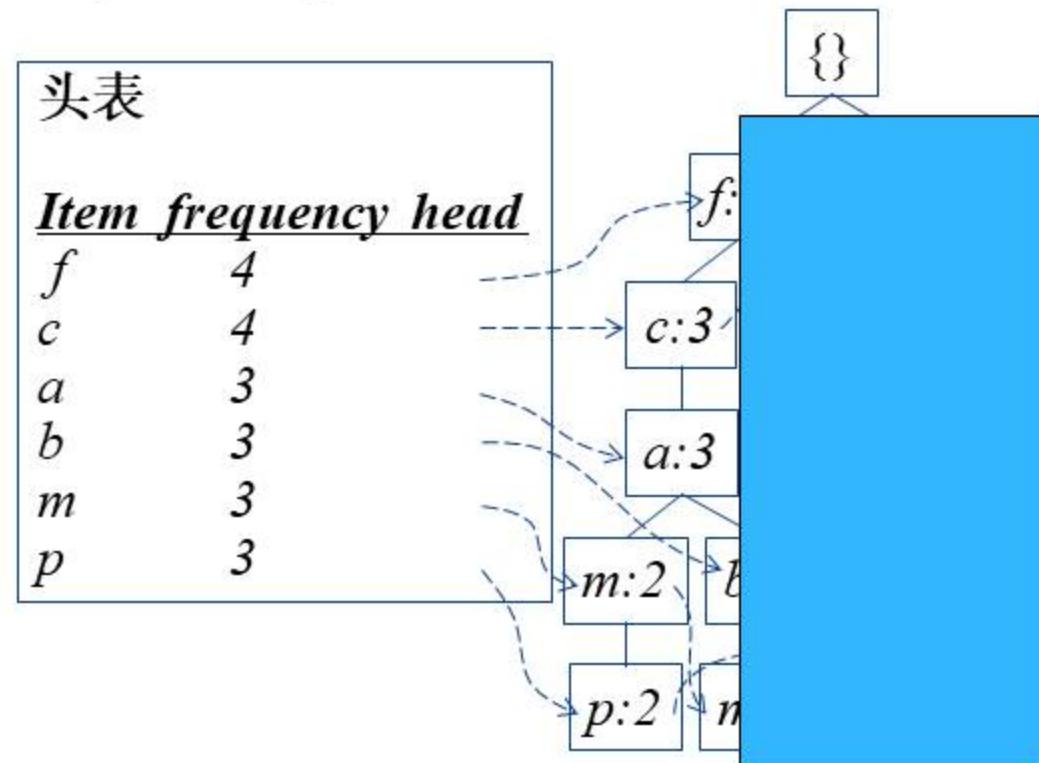
<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

最小支持度 = 0.5

步骤：

1. 扫描数据库一次，得到频繁1-项集
2. 把项按支持度递减排序
3. 再一次扫描数据库，建立FP-tree

头表	
<i>Item frequency head</i>	
f	4
c	4
a	3
b	3
m	3
p	3



练习：构造FP树

<u>TID</u>	<u>Items bought</u>	<u>(ordered) frequent items</u>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

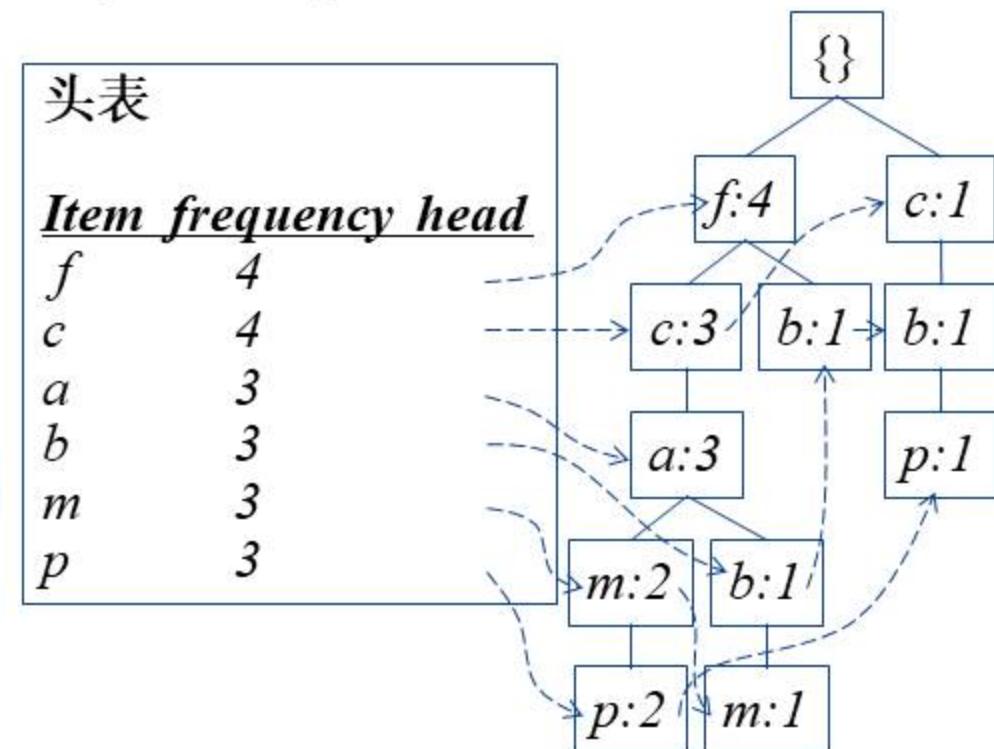
最小支持度 = 0.5

步骤：

1. 扫描数据库一次，得到频繁1-项集
2. 把项按支持度递减排序
3. 再一次扫描数据库，建立FP-tree

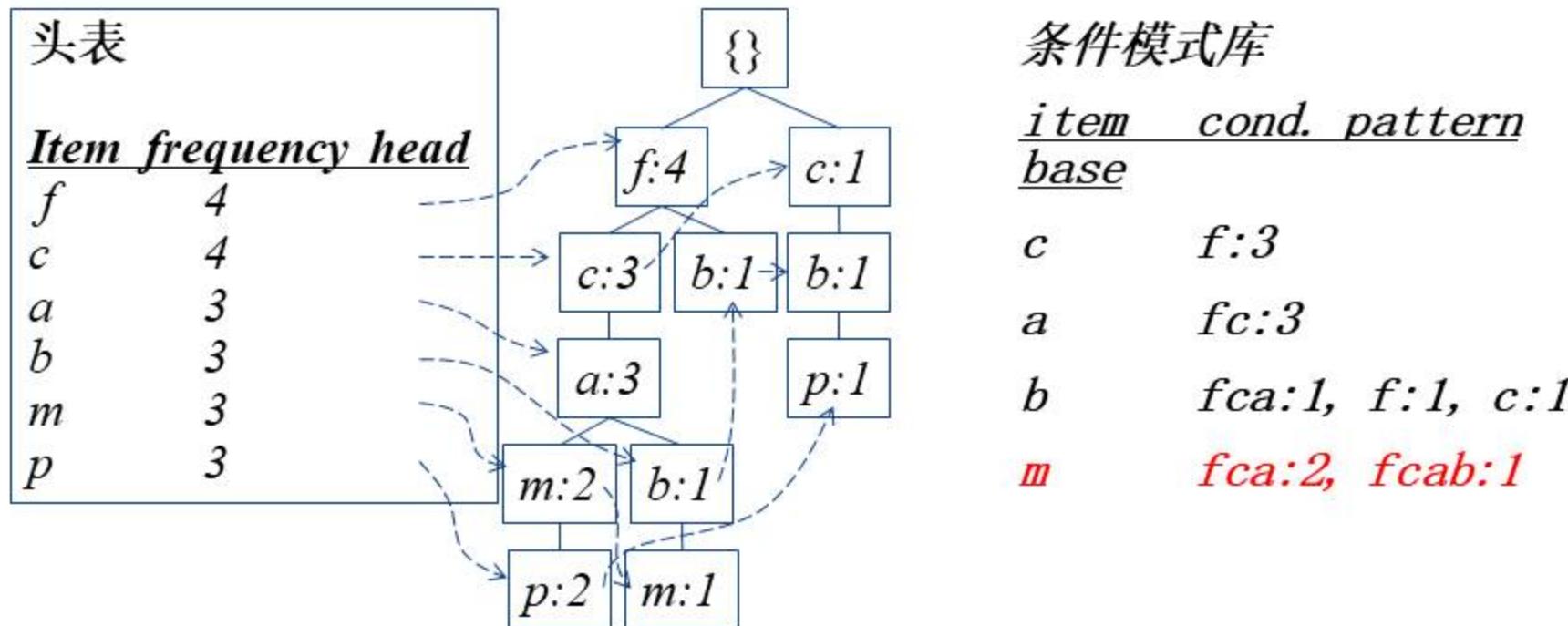
头表

<u>Item frequency head</u>	
f	4
c	4
a	3
b	3
m	3
p	3



练习：生成条件模式

- ◆ 从FP-tree的头表开始
- ◆ 按照每个频繁项的连接遍历 FP-tree
- ◆ 列出能够到达此项的所有前缀路径，得到条件模式基



P的条件模式基为：

A

f:2,c:2,a:2,m:2,p :2

B

f:2,c:1,a:1,m:1,p :1

C

f:1,c:1,a:1,m:1,p :1

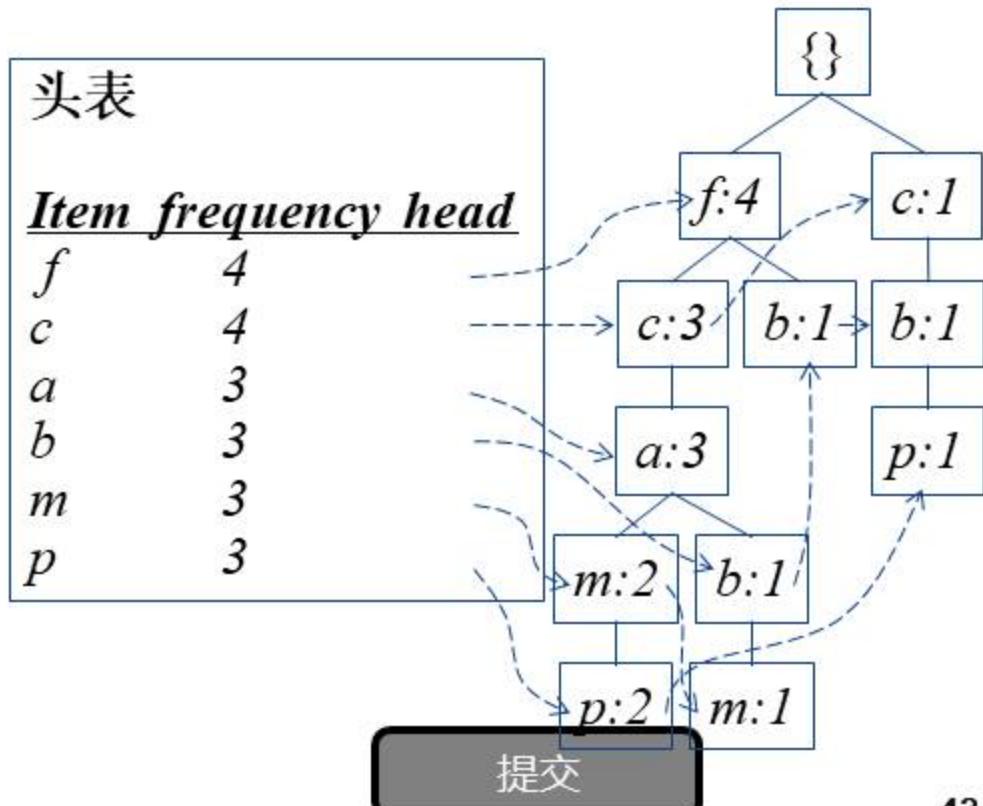
D

c:1,b:1,p :1

头表

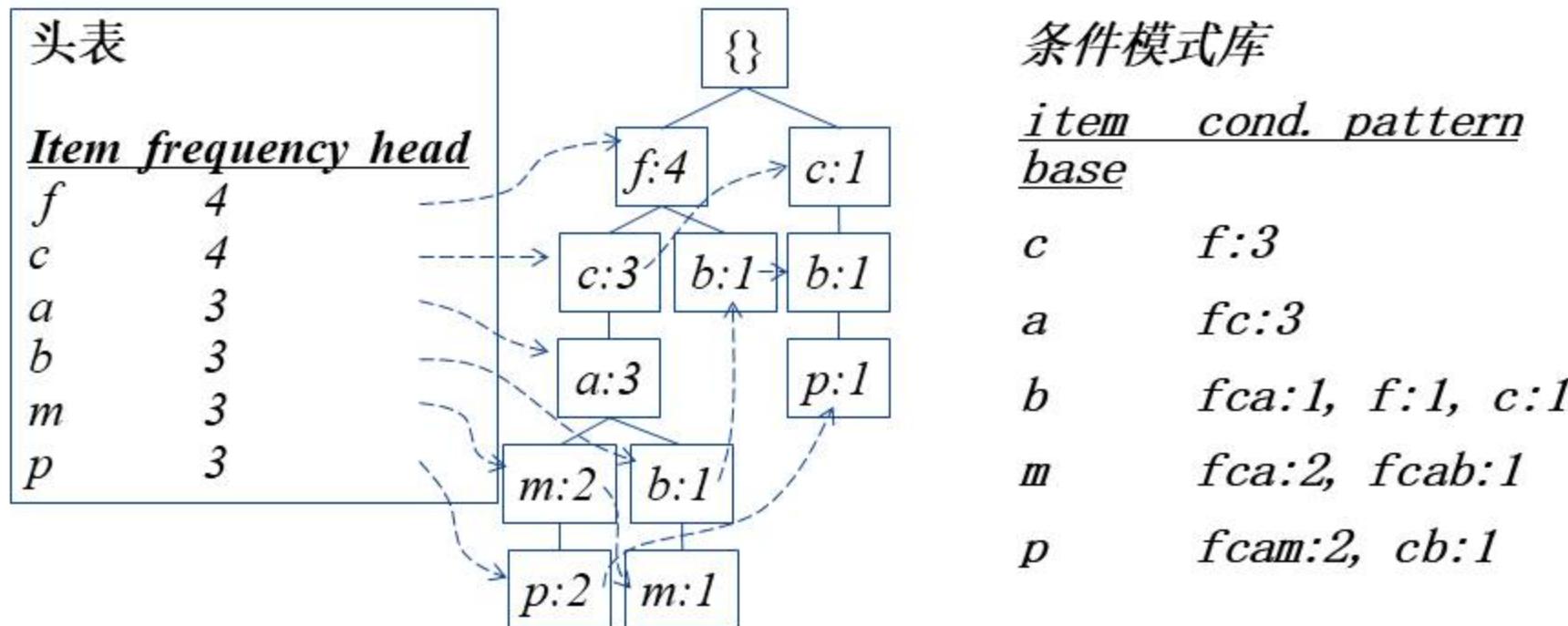
Item frequency head

f	4
c	4
a	3
b	3
m	3
p	3



练习：生成条件模式

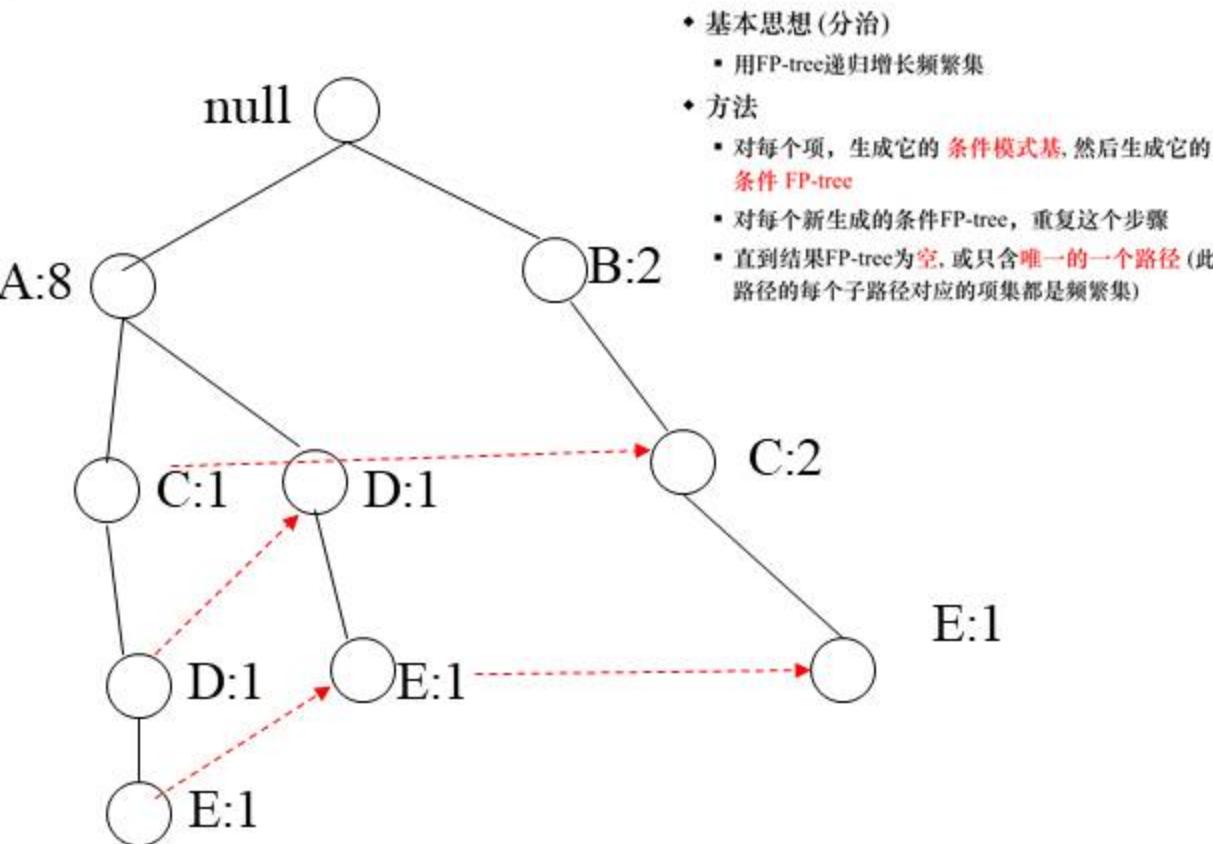
- ◆ 从FP-tree的头表开始
- ◆ 按照每个频繁项的连接遍历 FP-tree
- ◆ 列出能够到达此项的所有前缀路径，得到条件模式基



2.2构造FP树：用FP-tree挖掘频繁集

- ◆ 基本思想(分治)
 - 用FP-tree递归增长频繁集
- ◆ 方法
 - 对每个项，生成它的**条件模式基**，然后生成它的**条件 FP-tree**
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为**空**，或只含**唯一的一个路径**(此路径的每个子路径对应的项集都是频繁集)

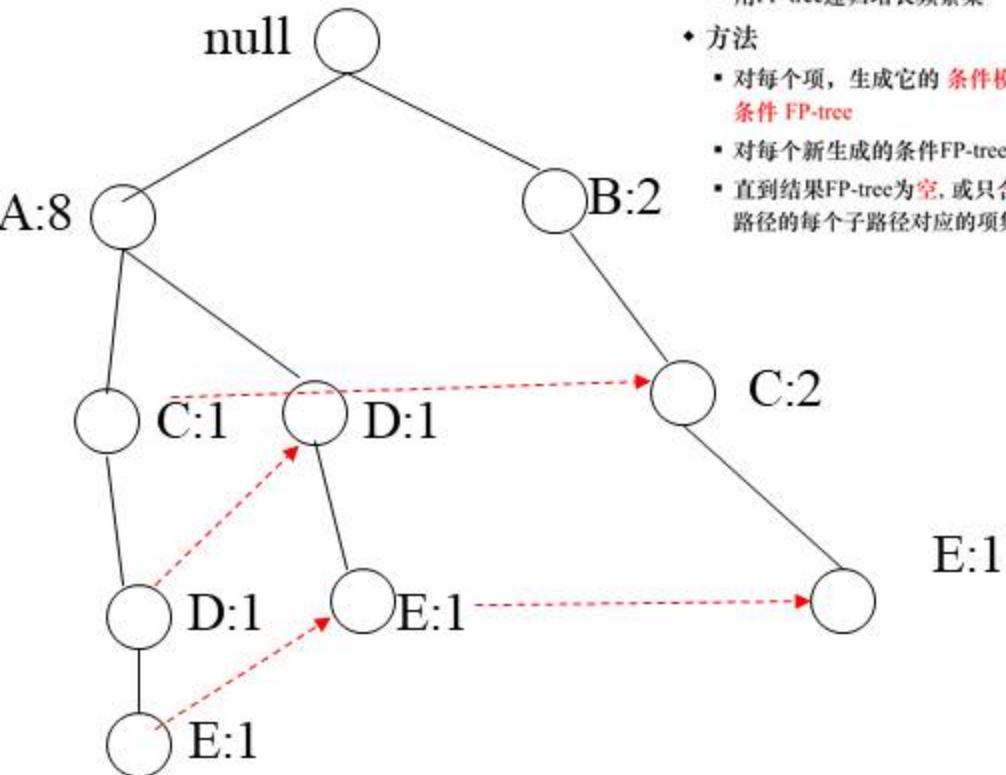
举例：节点E (1/3)



后缀模式
(PostModel)

E

举例：节点E (1/3)

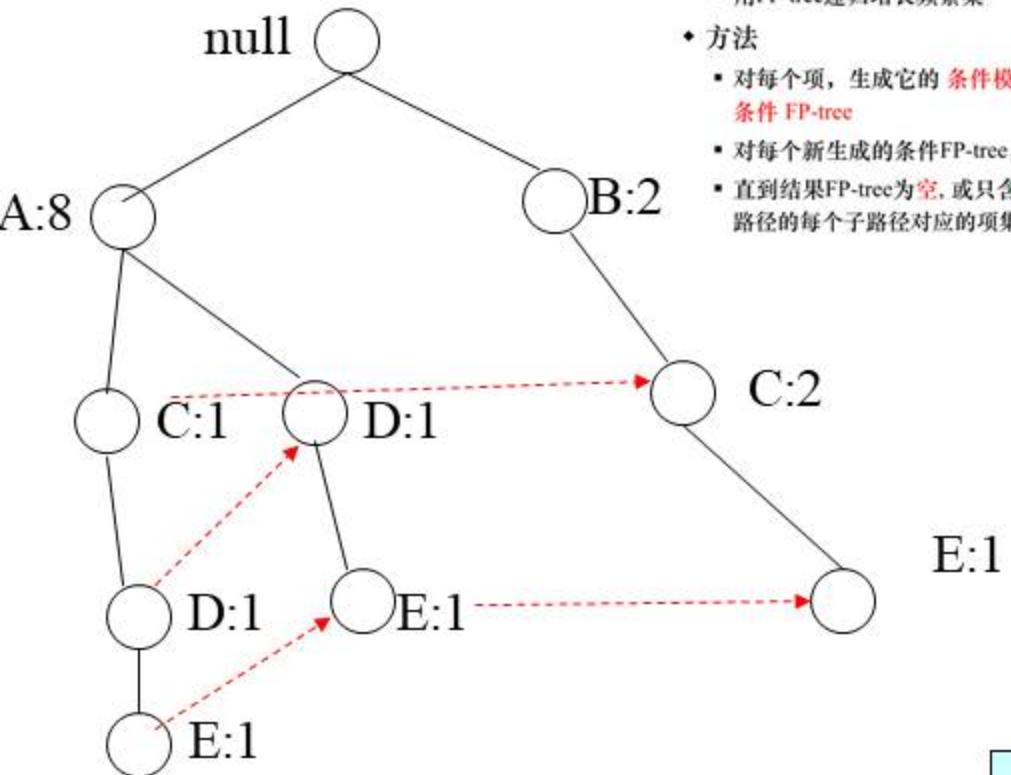


- 基本思想(分治)
 - 用FP-tree递归增长频繁集
- 方法
 - 对每个项，生成它的条件模式基，然后生成它的条件FP-tree
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为空，或只含唯一的一个路径(此路径的每个子路径对应的项集都是频繁集)

Items
{A:8,C:1,D:1,E:1}
{A:8,D:1,E:1}
{B:2,C:2,E:1}



举例：节点E (1/3)



- 基本思想(分治)
 - 用FP-tree递归增长频繁集
- 方法
 - 对每个项，生成它的**条件模式基**，然后生成它的**条件 FP-tree**
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为空，或只含**唯一的一个路径**(此路径的每个子路径对应的项集都是频繁集)

Items
{A:8,C:1,D:1,E:1}
{A:8,D:1,E:1}
{B:2,C:2,E:1}



Items
{A:1,C:1,D:1,E:1}
{A:1,D:1,E:1}
{B:1,C:1,E:1}



条件模式基 (Conditional Pattern Base,CPB)
{A:1,C:1,D:1}
{A:1,D:1}
{B:1,C:1}

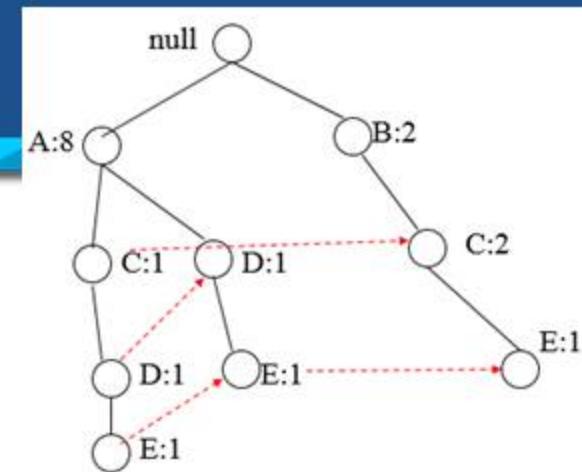
后缀模式 (PostModel)
E

举例：节点E (2/3)

- 基本思想(分治)
 - 用FP-tree递归增长频繁集
- 方法
 - 对每个项，生成它的条件模式基，然后生成它的条件 FP-tree
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为空，或只含唯一的一个路径(此路径的每个子路径对应的项集都是频繁集)

条件模式基
{A:1,C:1,D:1}
{A:1,D:1}
{B:1,C:1}

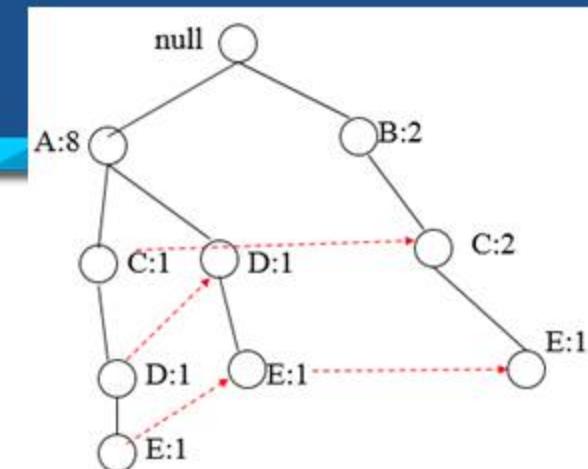
表头项
A:2
C:2
D:2



举例：节点E (2/3)

- 基本思想(分治)
 - 用FP-tree递归增长频繁集
- 方法
 - 对每个项，生成它的条件模式基，然后生成它的条件 FP-tree
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为空，或只含唯一的一个路径(此路径的每个子路径对应的项集都是频繁集)

条件模式基	
{A:1,C:1,D:1}	
	{A:1,D:1}
	{B:1,C:1}



表头项	
A:2	
C:2	
D:2	

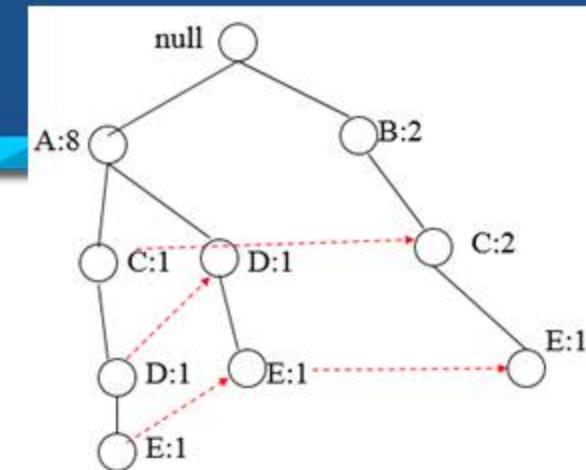
后缀模式 (PostModel)	
E	

频繁项集	
AE:2	
CE:2	
DE:2	
E:3	

举例：节点E (2/3)

- 基本思想(分治)
 - 用FP-tree递归增长频繁集
- 方法
 - 对每个项，生成它的 **条件模式基**，然后生成它的 **条件 FP-tree**
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为空，或只含**唯一的一个路径**(此路径的每个子路径对应的项集都是频繁集)

条件模式基	
{A:1,C:1,D:1}	
{A:1,D:1}	
{B:1,C:1}	



表头项	
A:2	
C:2	
D:2	

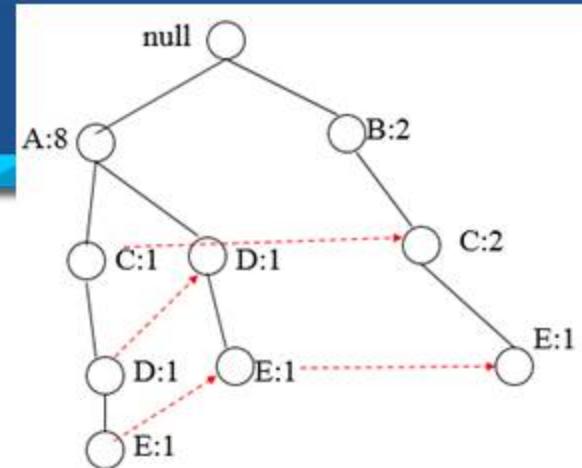
后缀模式 (PostModel)	
E	

频繁项集	
AE:2	
CE:2	
DE:2	
E:3	

新的后缀模式	
AE	
CE	
DE	

举例：节点E (3/3)

新的后缀模式
AE
CE
DE

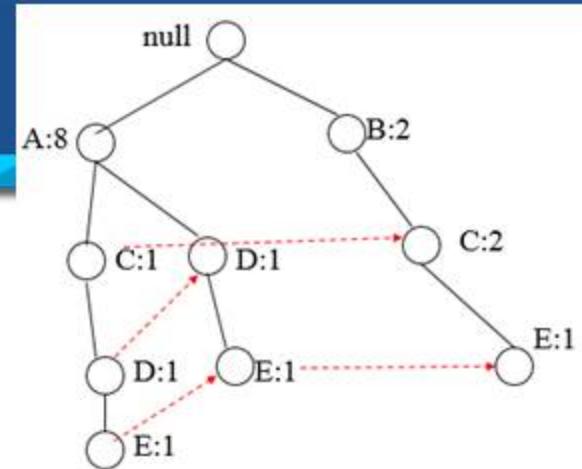


- 基本思想(分治)
 - 用FP-tree递归增长频繁集
- 方法
 - 对每个项，生成它的条件模式基，然后生成它的条件FP-tree
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为空，或只含唯一的一个路径(此路径的每个子路径对应的项集都是频繁集)

举例：节点E (3/3)

新的后缀模式
AE
CE
DE

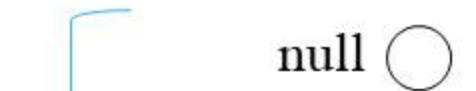
Items
{A:1,C:1,D:1,E:1}
{A:1,D:1,E:1}
{B:1,C:1,E:1}



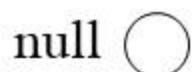
- 基本思想(分治)
 - 用FP-tree递归增长频繁集
- 方法
 - 对每个项，生成它的条件模式基，然后生成它的条件FP-tree
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为空，或只含唯一的一个路径(此路径的每个子路径对应的项集都是频繁集)

举例：节点E (3/3)

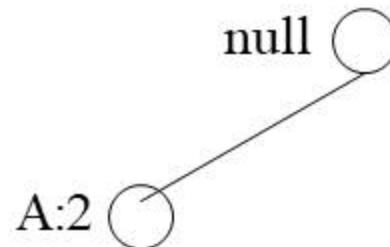
新的后缀模式
AE
CE
DE
Items
{A:1,C:1,D:1,E:1}
{A:1,D:1,E:1}
{B:1,C:1,E:1}



后缀模式为AE的FP-Tree

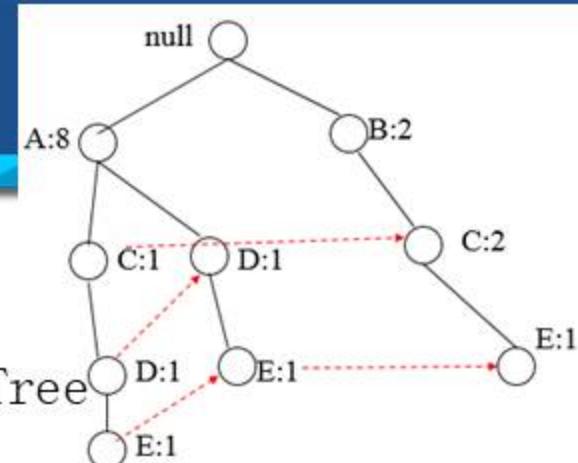


后缀模式为CE的FP-Tree



后缀模式为DE的FP-Tree

- 基本思想(分治)
 - 用FP-tree递归增长频繁集
- 方法
 - 对每个项，生成它的条件模式基，然后生成它的条件FP-tree
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为空，或只含唯一的一个路径(此路径的每个子路径对应的项集都是频繁集)

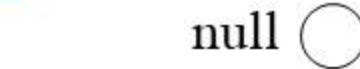


举例：节点E (3/3)

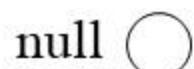
新的后缀模式
AE
CE
DE

Items
{A:1,C:1,D:1,E:1}
{A:1,D:1,E:1}
{B:1,C:1,E:1}

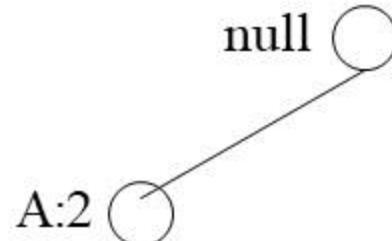
频繁项集
ADE:2
AE:2
CE:2
DE:2
E:3



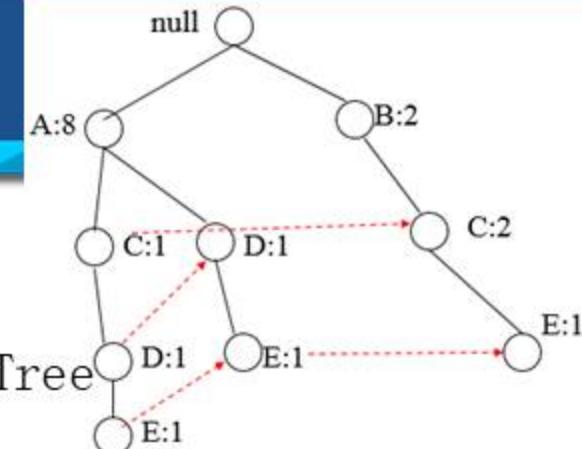
后缀模式为AE的FP-Tree



后缀模式为CE的FP-Tree



后缀模式为DE的FP-Tree

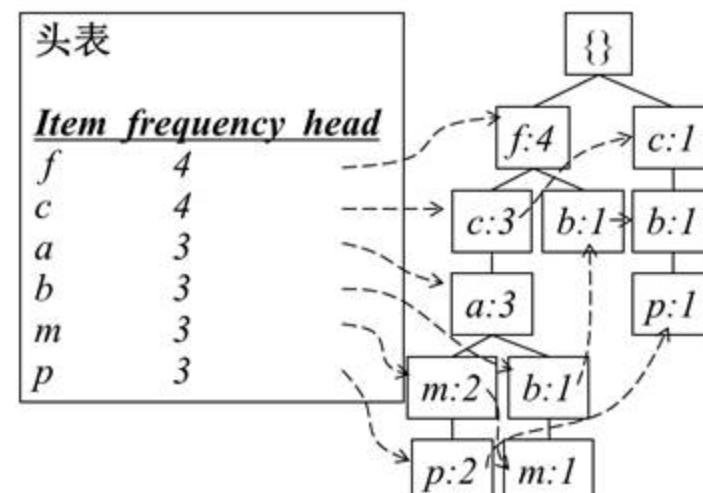


- 基本思想(分治)
 - 用FP-tree递归增长频繁集
- 方法
 - 对每个项，生成它的 **条件模式基**，然后生成它的 **条件 FP-tree**
 - 对每个新生成的条件FP-tree，重复这个步骤
 - 直到结果FP-tree为空，或只含**唯一的一个路径**(此路径的每个子路径对应的项集都是频繁集)

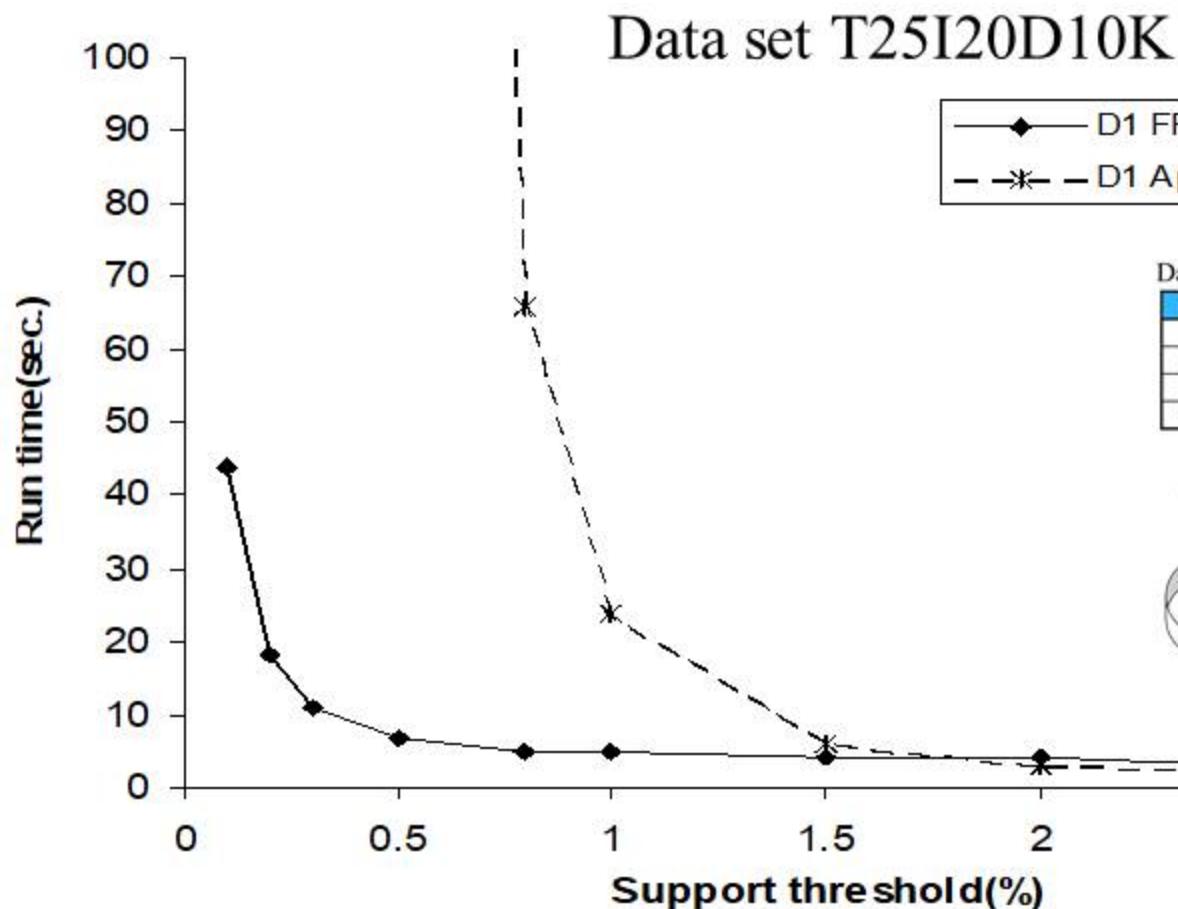
2.2FP-tree 结构的优点

- ◆ 完备:
 - 不会打破交易中的任何模式
 - 包含了频繁模式挖掘所需的全部信息
- ◆ 紧密
 - 支持度降序排列: 支持度高的项在FP-tree中共享的机会也高
 - 决不会比原数据库大

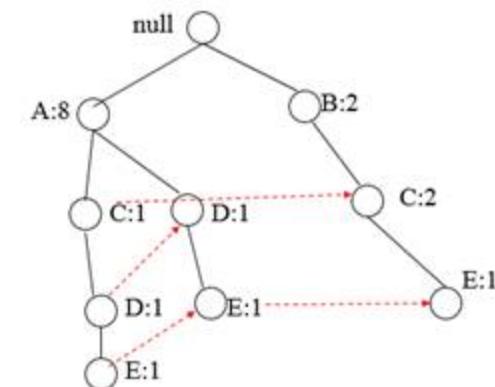
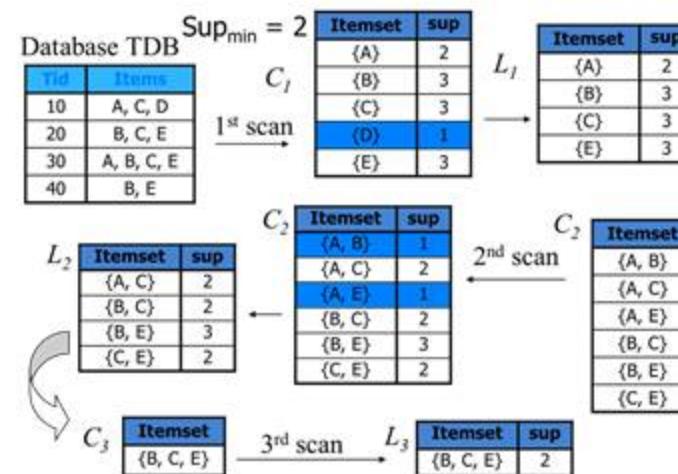
<https://www.cnblogs.com/itbuyixiaogong/p/9077428.html>



2.2FP-tree 结构的优点-性能对比



—◆— D1 FP-grow th runtime
—*— D1 Apriori runtime



2.3 挖掘关联规则 (Mining Association Rules)

- ◆ 大多数关联规则挖掘算法通常采用的一种策略是，将关联规则挖掘任务分解为如下两个主要的子任务：
 1. 频繁项集产生 (Frequent Itemset Generation)
 - 其目标是发现满足最小支持度阈值的所有项集，这些项集称作频繁项集。
 2. 规则的产生 (Rule Generation)
 - 其目标是从上一步发现的频繁项集中提取所有高置信度的规则，这些规则称作强规则 (strong rule)。

2.3产生关联规则

- ◆ 任务描述：给定频繁项集 Y ，查找 Y 的所有非空真子集 $X \subset Y$ ，使得 $X \rightarrow Y - X$ 的置信度超过最小置信度阈值 $minconf$
 - 例子：If $\{A,B,C\}$ is a frequent itemset, 候选规则如下：
$$\begin{array}{lll} AB \rightarrow C, & AC \rightarrow B, & BC \rightarrow A \\ A \rightarrow BC, & B \rightarrow AC, & C \rightarrow AB \end{array}$$

- ◆ 如果 $|Y| = k$, 那么会有 $2^k - 2$ 个候选关联规则
(不包括 $Y \rightarrow \emptyset$ and $\emptyset \rightarrow Y$)

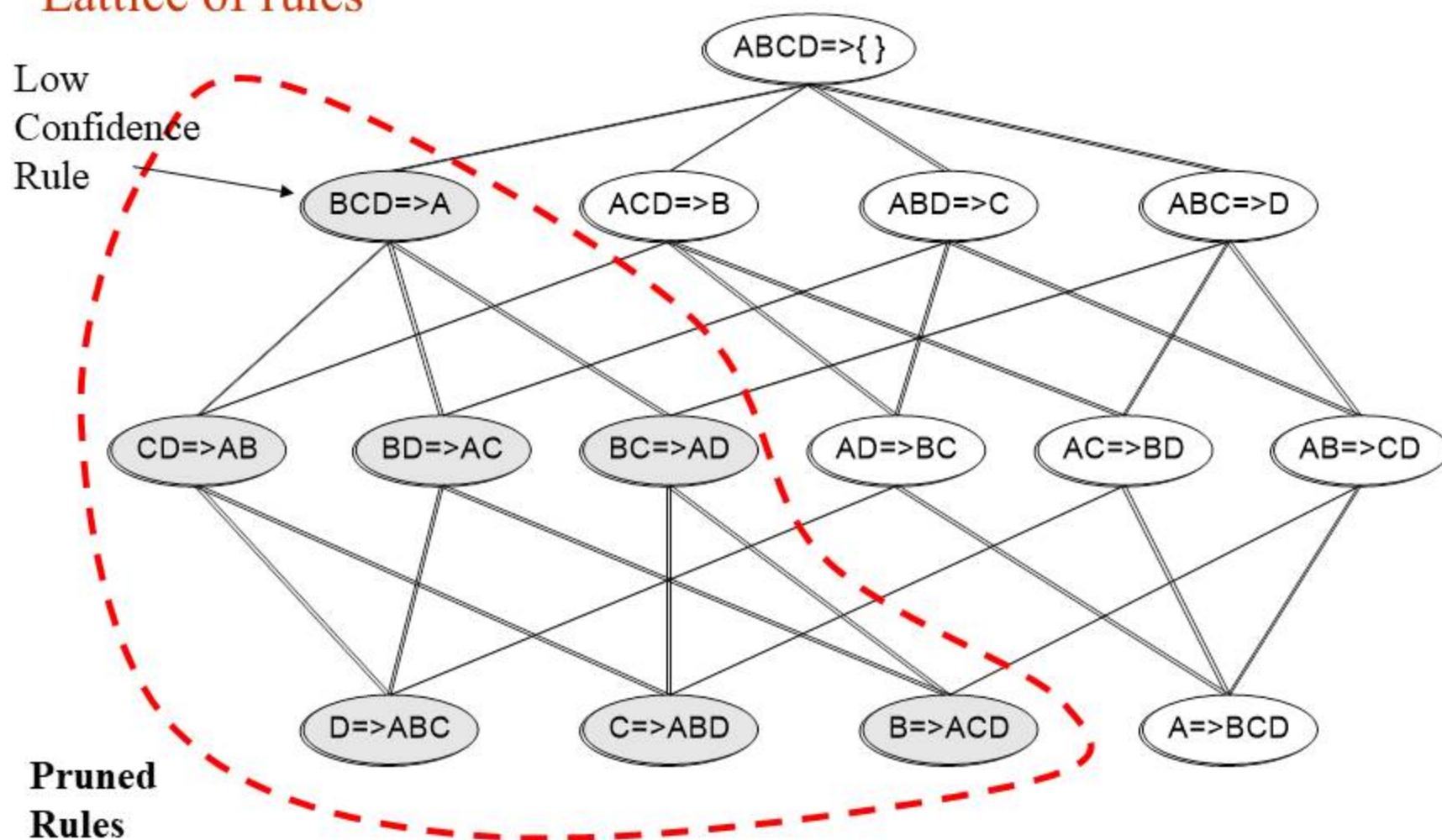
2.3产生关联规则

- ◆ How to efficiently generate rules from frequent itemsets?
 - 通常，置信度不满足反单调性 (anti-monotone property)，例如： $c(ABC \rightarrow D)$ 可能大于也可能小于 $c(AB \rightarrow D)$
 - 但是，针对同一个频繁项集的关联规则，如果规则的后件满足子集关系，那么这些规则的置信度间满足反单调性
 - e.g., $Y = \{A, B, C, D\}$:
 $c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$

2.3 Rule Generation for Apriori Algorithm

Lattice of rules

Low
Confidence
Rule



主要内容

- ◆ 1. 基于概念
- ◆ 2. 频繁项挖掘算法
- ◆ 3. 关联分析的评估

3. 关联分析的评估 (Pattern Evaluation)

- $play\ basketball \Rightarrow eat\ cereal$ [支持度=, 置信度=]

	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

3. 关联分析的评估 (Pattern Evaluation)

- $play\ basketball \Rightarrow eat\ cereal [40%, 66.7%]$ is misleading
 - The overall % of students eating cereal is $75\% > 66.7\%$.
- $play\ basketball \Rightarrow not\ eat\ cereal [20%, 33.3\%]$ is more accurate, although with lower support and confidence
- Measure of dependent/correlated events: lift

$$lift = \frac{P(A \cup B)}{P(A)P(B)}$$

$$lift(B, C) = \frac{2000/5000}{3000/5000 * 3750/5000} = 0.89$$

	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

$$lift(B, \neg C) = \frac{1000/5000}{3000/5000 * 1250/5000} = 1.33$$

- ◆ 1. 基于概念
- ◆ 2. 频繁项挖掘算法
- ◆ 3. 关联分析的评估

问题？

第5次课后作业

- ◆ 第五次课后作业-在educoder平台上完成作业
- ◆ <https://www.educoder.net/shixuns/egbxla2s/challenges>
- ◆ <https://www.educoder.net/shixuns/7ctykufv/challenges>

提交作业截至时间：2020年3月05日

zyding1983 / datamining

 Code

 Issues 0

 Pull requests 0

 Actions

 Projects 0

Branch: master ▾

datamining / 2020A / codes / chapters4apirio /



zyding1983 add dm ppt all pdf

..

 apriori-python-master

add ne1'

 .DS_Store

add dm

 apriori.py

add ne1'

```
▶ 69
▶ 70         U = []
▶ 71     for f in F:
▶ 72         for x in f:
▶ 73             U.append(x)
▶ 74     return U
▶ 75
▶ 76
▶ 77 T = [['A', 'C', 'D'], ['B', 'C', 'E'], ['A', 'B', 'C', 'E'], ['B', 'E']]
▶ 78
▶ 79 Z= apriori(T,3)
▶ 80 print(Z)
```

```
apriori() > for f in D[0] > if C[f]>=minSupport
```

un: apriori

```
/Users/zyding/venv/bin/python /Users/zyding/IdeaProjects/python1/untitled/apr
[['B'], ['C'], ['E'], ['B', 'E']]
```

```
Process finished with exit code 0
```

```
▶ 60     F.append([])
▶ 61     for c in D[k]:
▶ 62         count = 0;
▶ 63         for t in T:
▶ 64             if compareList(c,t):
▶ 65                 count += 1
▶ 66             if count >= minSupport:
▶ 67                 F[k].append(c)
▶ 68             k += 1
▶ 69
▶ 70     U = []
▶ 71     for f in F:
▶ 72         for x in f:
▶ 73             U.append(x)
▶ 74     return U
▶ 75
▶ 76
▶ 77 T = [['A', 'C', 'D'], ['B', 'C', 'E'], ['A', 'B', 'C', 'E'], ['B', 'E']]
▶ 78
▶ 79 Z= apriori(T,2)
▶ 80 print(Z)
```

apriori()

Run: apriori

```
/Users/zyding/venv/bin/python /Users/zyding/IdeaProjects/python1/untitled/apriori.py
[['A'], ['B'], ['C'], ['E'], ['A', 'C'], ['B', 'C'], ['B', 'E'], ['C', 'E'], ['B', 'C', 'E']]
```

zyding1983 / datamining

 Code

 Issues 0

 Pull requests 0

 Actions

 Projects 0

Branch: master ▾

[datamining / 2020A / codes / chapters4apirio /](#)



zyding1983 add dm ppt all pdf

..



apriori-python-master

add ne1'



.DS_Store

add dm



apriori.py

add ne1'

- apriori-python-master
 - data
 - 1000
 - 5000
 - 20000
 - 75000
 - example
 - transcription
- docs
 - apriori.py
 - boxF.py
 - goods.csv
 - goods-BK.csv
 - lab3.466.pdf
 - out1.csv
 - pcaDemo2.py
- README.md
- aprioridemo.py

```
B I M View
1 Association Rules Mining, Apriori Implementation
2
3 Timothy Asp, Caleb Carlton
4
5 Input Format: python apriori.py [--no-rules] <dataFile-out1.csv> <minSup> <minConf>
6 --no-rules will run the code without rules generation.
7 The input datafile must be in the sparse vector format (see *-out1.csv in the different
8
9 Example:
10
11 python apriori.py --no-rules data/1000/1000-out1.csv .03 .7
12 python apriori.py data/example/out1.csv .03 .7
13
14 == Extended Bakery Printouts
15
16 =====
17 Dataset: data/example/out1.csv MinSup: 0.03 MinConf: 0.7
18 =====
19 1 : Blackberry Tart (15), Apple Danish (36) support= 0.139
20 2 : Gongolais Cookie (22), Napoleon Cake (9) support= 0.181
21 3 : Lemon Cake (1), Single Espresso (49) support= 0.127
22 4 : Apple Tart (12), Berry Tart (14), Blueberry Tart (16) support= 0.257
23
```

```
▼ apriori-python-master
  ► data
  ► docs
  ► apriori.py
  ► boxF.py
  ► goods.csv
  ► goods-BK.csv
  ► lab3.466.pdf
  ► out1.csv
  ► pcaDemo2.py
  ► README.md
  ► aprioridemo.py
```

```
► yingyong111
```

```
apriori (1) ×
```

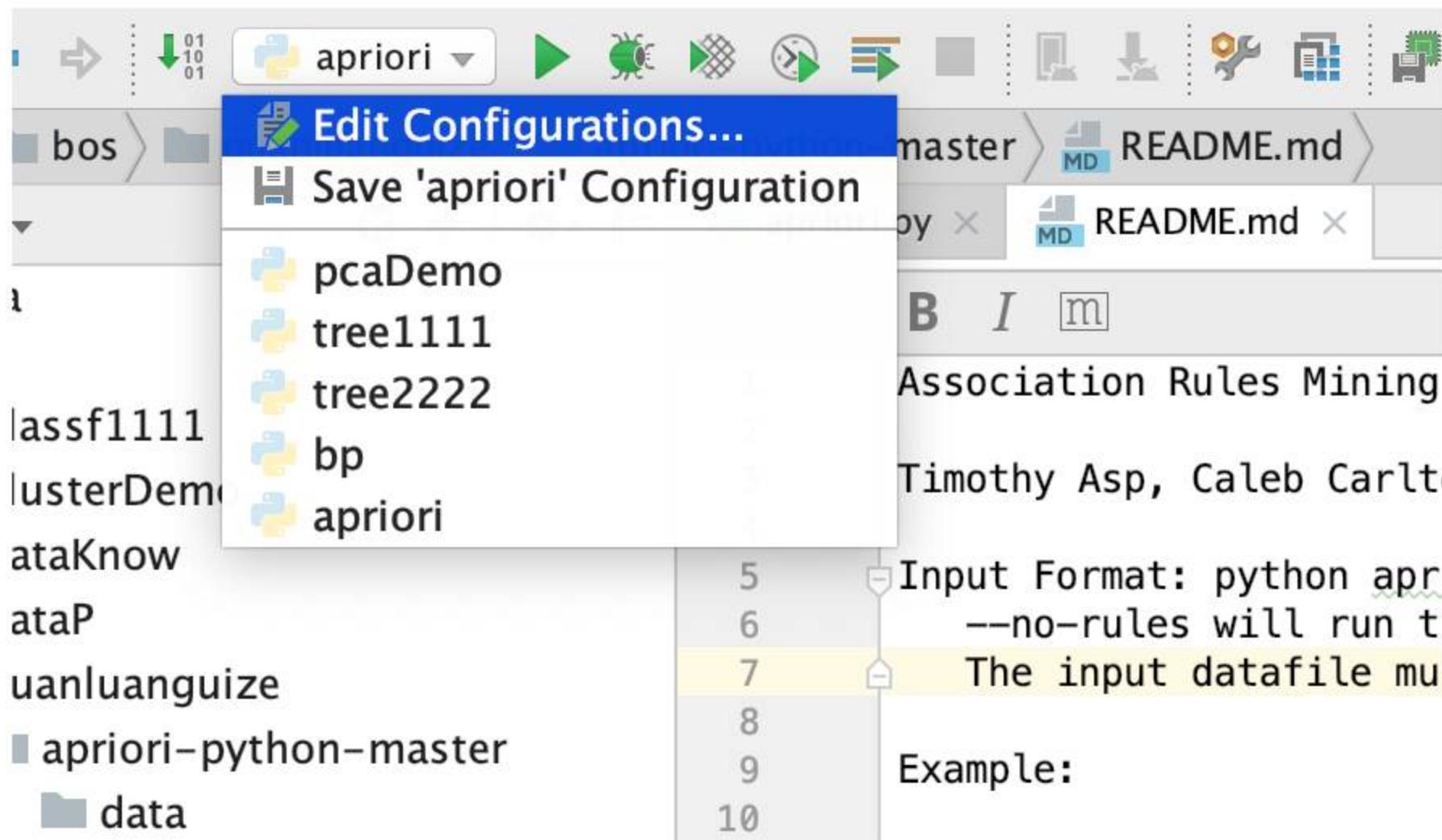
```
/Users/zyding/venv/bin/python /Users/zyding/IdeaProjects/python1/bos/guanluangu  
Expected input format: python apriori.py <dataset.csv> <minSup> <minConf>
```

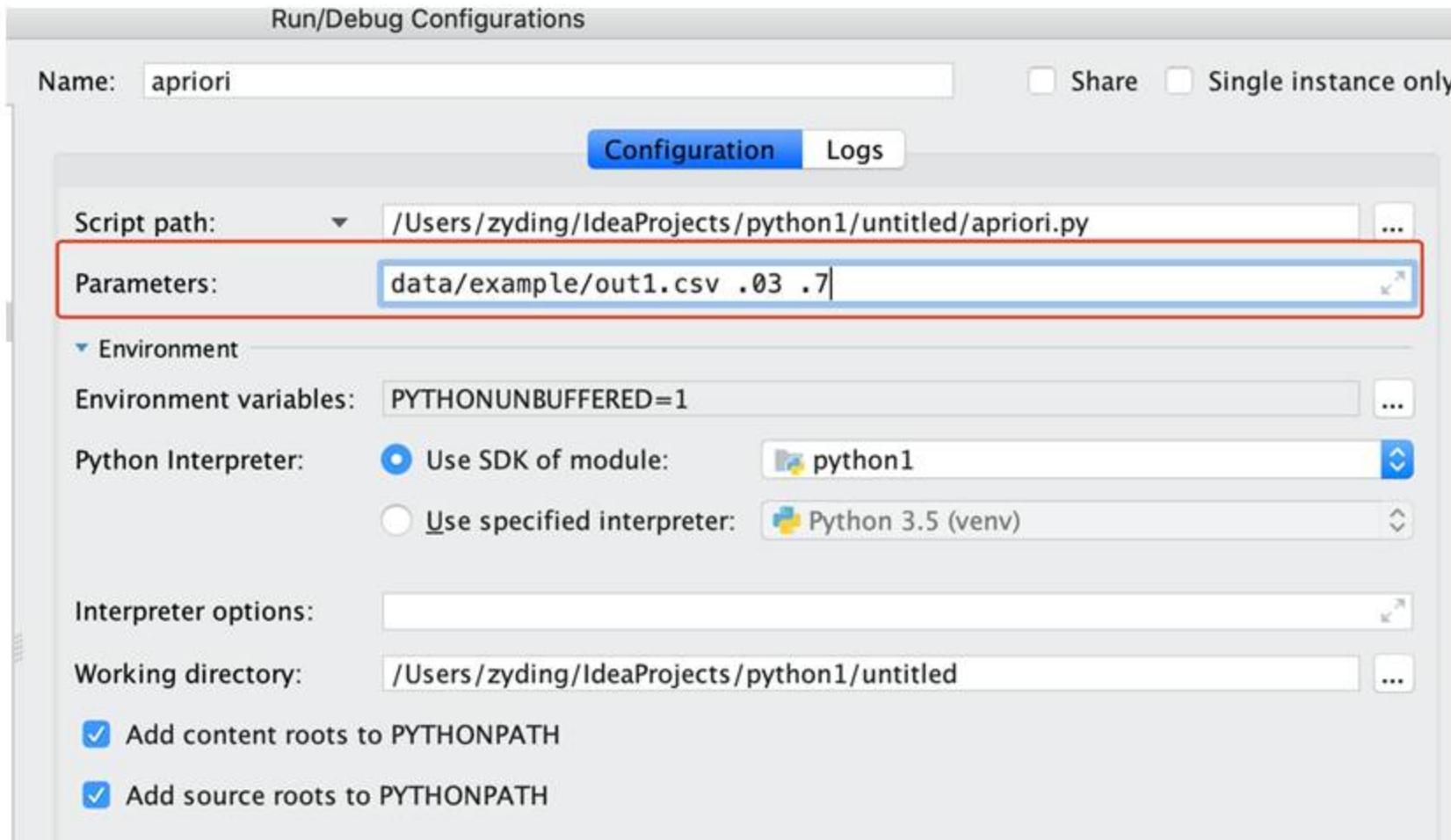
```
Process finished with exit code 0
```

```
    11
    12
    13
    14
    15
    16
    17
    18
    19
    20
    21
    22
    23
    24
    25
    26
    27
```

```
self.udata = udata
self.transList = defaultdict(list)
self.freqList = defaultdict(int)
self.itemset = set()
self.highSupportList = list()
self.numItems = 0
self.prepData() # init
self.F = defaultdict(list)
self.minSup = minSup
self.minConf = minConf

def genAssociations(self):
    candidate = {}
    count = {}
```

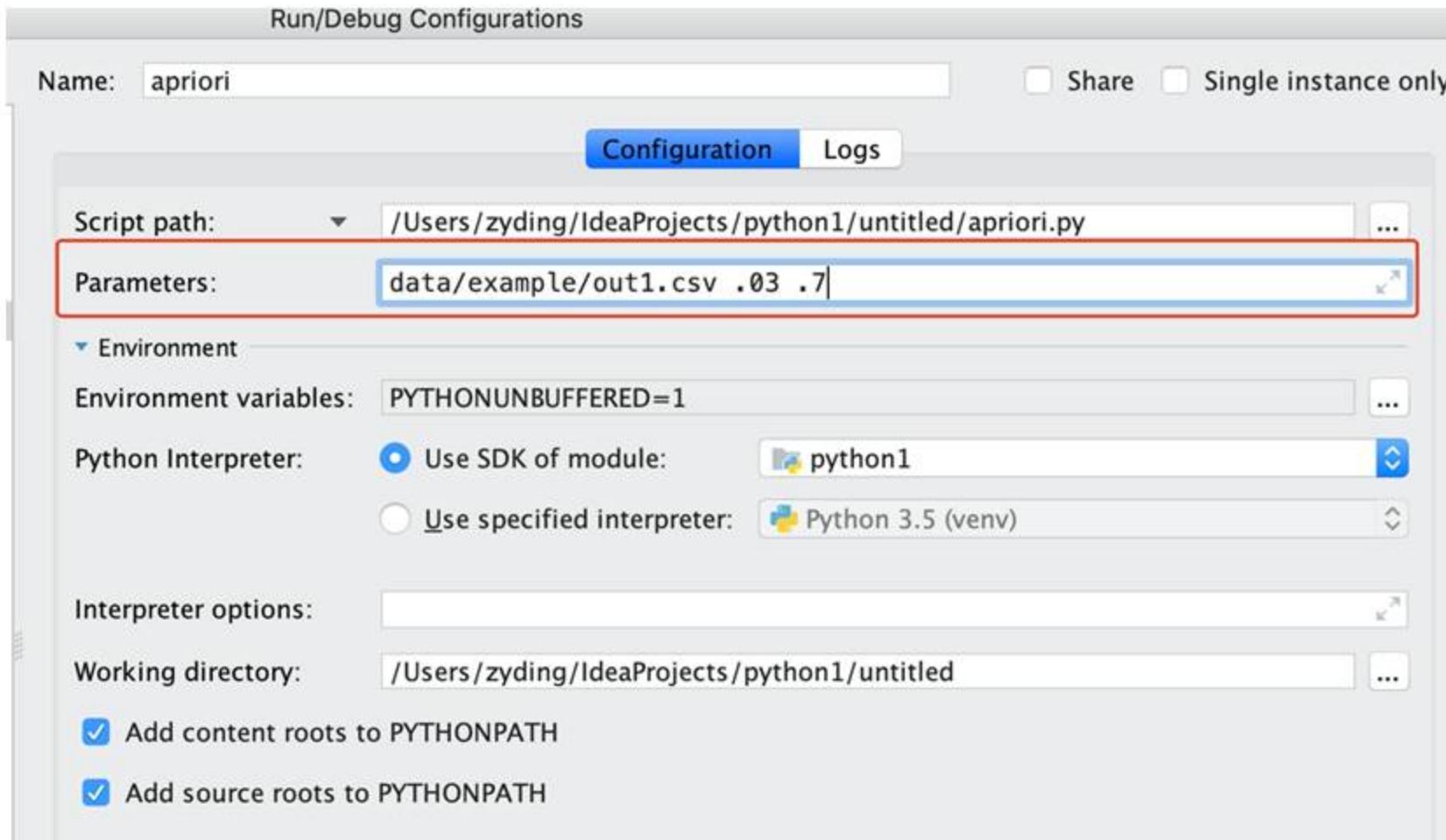




Run: apriori (1) ×

```
/Users/zyding/venv/bin/python /Users/zyding/IdeaProjects/python1/bos/guanluanguize/apriori-python-master/apriori.py data/
Dataset: data/example/out1.csv MinSup: 0.03 MinConf: 0.7
=====
1 : Lemon 0000 (1), Single KKKK (49) support= 0.127
2 : Gongolais HHHH (22), Napoleon 0000 (9) support= 0.181
3 : Apple GGGG (12), Berry GGGG (14) support= 0.268
4 : Blackberry GGGG (15), Apple RRRR (36) support= 0.139
5 : Apple GGGG (12), Blueberry GGGG (16) support= 0.259
6 : Berry GGGG (14), Blueberry GGGG (16), Apple GGGG (12) support= 0.257
Skyline Itemsets: 6
Rule 1 : Lemon 0000 (1) --> Single KKKK (49) [sup= 0.127 conf= 0.8141025641025641 ]
Rule 2 : Single KKKK (49) --> Lemon 0000 (1) [sup= 0.127 conf= 0.7839506172839507 ]
Rule 3 : Gongolais HHHH (22) --> Napoleon 0000 (9) [sup= 0.181 conf= 0.8418604651162791 ]
Rule 4 : Napoleon 0000 (9) --> Gongolais HHHH (22) [sup= 0.181 conf= 0.8044444444444444 ]
Rule 5 : Apple GGGG (12) --> Berry GGGG (14) [sup= 0.268 conf= 0.8933333333333333 ]
Rule 6 : Berry GGGG (14) --> Apple GGGG (12) [sup= 0.268 conf= 0.9146757679180887 ]
Rule 7 : Blackberry GGGG (15) --> Apple RRRR (36) [sup= 0.139 conf= 0.7513513513513513 ]
Rule 8 : Apple RRRR (36) --> Blackberry GGGG (15) [sup= 0.139 conf= 0.7988505747126436 ]
Rule 9 : Apple GGGG (12) --> Blueberry GGGG (16) [sup= 0.259 conf= 0.8633333333333333 ]
Rule 10 : Blueberry GGGG (16) --> Apple GGGG (12) [sup= 0.259 conf= 0.9151943462897526 ]
Rule 11 : Berry GGGG (14), Blueberry GGGG (16) --> Apple GGGG (12) [sup= 0.257 conf= 0.9961240310077519 ]
```

Process finished with exit code 0



Run/Debug Configurations

Name: apriori Share Single instance only

Configuration Logs

Script path: /Users/zyding/IdeaProjects/python1/untitled/apriori.py ...

Parameters: data/example/out1.csv .03 .7

Environment

apriori-python-master

- ▶ data
- ▶ docs
- ▶ apriori.py
- ▶ boxF.py
- ▶ goods.csv
- ▶ goods-BK.csv
- ▶ lab3.466.pdf
- ▶ out1.csv
- ▶ pcaDemo2.py
- ▶ README.md
- ▶ aprioridemo.py

Plugins supporting *.csv files found.

1	1, 2, 5, 27, 29, 32, 48
2	2, 15, 27, 31, 33, 36
3	3, 9, 16, 22
4	4, 12, 14, 16
5	5, 2, 9, 44
6	6, 9, 26, 39
7	7, 10, 31
8	8, 1, 47, 49
9	9, 12, 14, 16
10	10, 12, 14, 16
11	11, 1, 19, 25, 49
12	12, 15, 23, 36
13	13, 15, 17, 36
14	14, 15, 18, 36, 37, 41

▼	apriori-python-master
▼	data
►	1000
►	5000
►	20000
►	75000
►	example
►	transcription
►	docs
	apriori.py
	boxF.py
	goods.csv
	goods-BK.csv
	lab3.466.pdf
	out1.csv
	pcaDemo2.py
	README.md
	aprioridemo.py

Plugins supporting *.csv files found.

```

1  0,'Chocolate','0000',8.95,'AAA'
2  1,Lemon','0000',8.95,'AAA'
3  2,'Casino','0000',15.95,'AAA'
4  3,'Opera','0000',15.95,'AAA'
5  4,'Strawberry','0000',11.95,'AAA'
6  5,'Truffle','0000',15.95,'AAA'
7  6,'Chocolate','FFFF',3.25,'AAA'
8  7,'GGGG','FFFF',3.5,'AAA'
9  8,'Vanilla','FFFF',3.25,'AAA'
10 9,'Napoleon','0000',13.49,'AAA'
11 10,'Almond','GGGG',3.75,'AAA'
12 11,'Apple','NNNN',5.25,'AAA'
13 12,'Apple','GGGG',3.25,'AAA'
14 13,'Apricot','GGGG',3.25,'AAA'
15 14,'Berry','GGGG',3.25,'AAA'
16 15,'Blackberry','GGGG',3.25,'AAA'
17 16,'Blueberry','GGGG',3.25,'AAA'
18 17,'Chocolate','GGGG',3.75,'AAA'
19 18,'Cherry','GGGG',3.25,'AAA'
20 19,'Lemon','GGGG',3.25,'AAA'
21 20,'Pecan','GGGG',3.75,'AAA'
22 21,'Ganache','HHHH',1.15,'AAA'
...

```

Run: apriori (1) ×

```
/Users/zyding/venv/bin/python /Users/zyding/IdeaProjects/python1/bos/guanluanguize/apriori-python-master/apriori.py data/
Dataset: data/example/out1.csv MinSup: 0.03 MinConf: 0.7
=====
1 : Lemon 0000 (1), Single KKKK (49) support= 0.127
2 : Gongolais HHHH (22), Napoleon 0000 (9) support= 0.181
3 : Apple GGGG (12), Berry GGGG (14) support= 0.268
4 : Blackberry GGGG (15), Apple RRRR (36) support= 0.139
5 : Apple GGGG (12), Blueberry GGGG (16) support= 0.259
6 : Berry GGGG (14), Blueberry GGGG (16), Apple GGGG (12) support= 0.257
Skyline Itemsets: 6
Rule 1 : Lemon 0000 (1) --> Single KKKK (49) [sup= 0.127 conf= 0.8141025641025641 ]
Rule 2 : Single KKKK (49) --> Lemon 0000 (1) [sup= 0.127 conf= 0.7839506172839507 ]
Rule 3 : Gongolais HHHH (22) --> Napoleon 0000 (9) [sup= 0.181 conf= 0.8418604651162791 ]
Rule 4 : Napoleon 0000 (9) --> Gongolais HHHH (22) [sup= 0.181 conf= 0.8044444444444444 ]
Rule 5 : Apple GGGG (12) --> Berry GGGG (14) [sup= 0.268 conf= 0.8933333333333333 ]
Rule 6 : Berry GGGG (14) --> Apple GGGG (12) [sup= 0.268 conf= 0.9146757679180887 ]
Rule 7 : Blackberry GGGG (15) --> Apple RRRR (36) [sup= 0.139 conf= 0.7513513513513513 ]
Rule 8 : Apple RRRR (36) --> Blackberry GGGG (15) [sup= 0.139 conf= 0.7988505747126436 ]
Rule 9 : Apple GGGG (12) --> Blueberry GGGG (16) [sup= 0.259 conf= 0.8633333333333333 ]
Rule 10 : Blueberry GGGG (16) --> Apple GGGG (12) [sup= 0.259 conf= 0.9151943462897526 ]
Rule 11 : Berry GGGG (14), Blueberry GGGG (16) --> Apple GGGG (12) [sup= 0.257 conf= 0.9961240310077519 ]
```

Process finished with exit code 0