



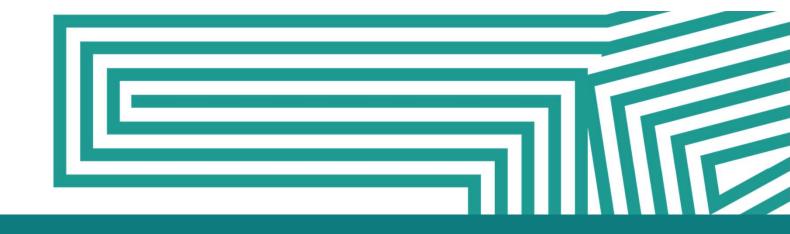


### 目录

- 一 任务背景
- 二任务描述
  - 1. 基本任务
  - 2. 进阶任务
- 三算法流程
- 四 验收流程







## 一任务背景



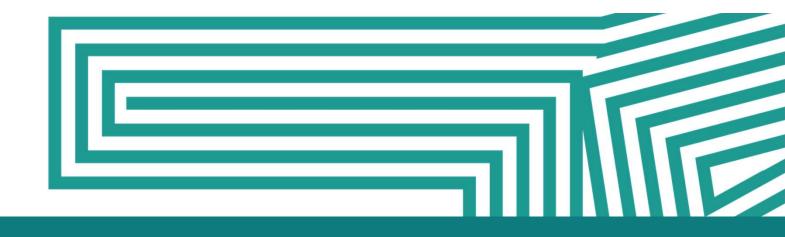
### 一 任务背景



- 1、理解Apriori算法思想与流程;
- 2、应用Apriori思想解决问题;

5/22/2024 4





## 二任务描述



### 二任务描述



#### ◆基本任务

- 以实验1得到的频率前1000个title及其相互引用关系作为关系挖掘实验的数据,从实验一中得到的引用关系数据为<<title,<title1,...,titlek>>,...>,将其处理为<<title,title1,...,titlek>,...>作为算法输入。
- 编程实现Apriori算法,使用实验数据进行实验,获得频繁项集以及关联规则。
- 输出1~4阶频繁项集与关联规则,各个频繁项的支持度,各个规则的置信度,各阶频繁项集的数量以及关联规则的总数,固定参数以方便检查,频繁项集的最小支持度为0.15,关联规则的最小置信度为0.3。此处支持度的定义为某个项集出现的频率,也就是包含该项集的数目与总数目的比例(总的购物篮数目为1000)。

5/22/2024 6

### 二任务描述



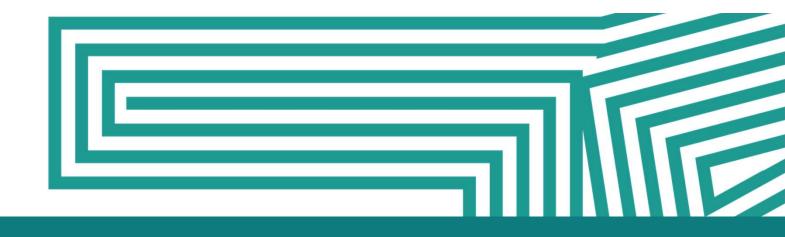
#### ◆进阶任务

• 在Apriori算法的基础上,使用pcy算法对二阶频繁项集的计算阶段进行优化。

• 除基本任务的输出外,额外输出pcy算法中的vector的值,以bit位的形式输出。固定参数以方便检查,频繁项集的最小支持度为0.15,关联规则的最小置信度为0.3。

5/22/2024 7





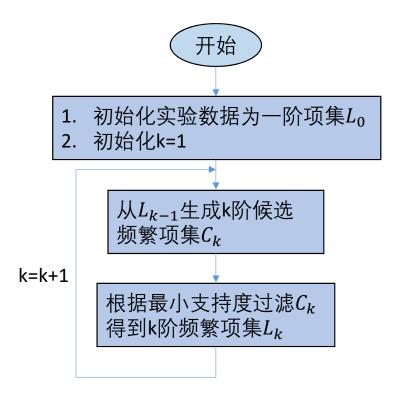
## 三算法流程



### 三 算法流程



◆挖掘频繁项集



### 三 算法流程



#### ◆由频繁项集产生关联规则

- 对于一个频繁项I, I的每个子集A, 生成一个规则 $A \rightarrow I \setminus A$
- 筛选出所有置信度大于最小置信度的规则
- 一个规则的置信度计算公式:  $confidence(A \rightarrow I \setminus A) = support(I) / support(A)$

### 三 算法流程



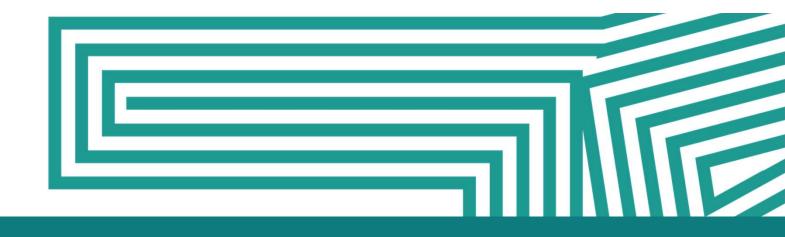
◆进阶: PCY

```
Frequent items
                                                        Item counts
 FOR (each basket):
                                                                         Bitmap
                                                   Main memory
      FOR (each item in the basket) :
                                                        Hash table
                                                                        Counts of
                                                         for pairs
                                                                        candidate
            add 1 to item's count;
                                                                          pairs
     FOR (each pair of items) :
New
            hash the pair to a bucket;
                                                        Pass 1
                                                                        Pass 2
PCY
            add 1 to the count for that bucket;
```

bit vector 的每一位代表一个bucket是否为频繁的,**如果一个bucket中的计数小于最小 支持度,那么映射到这个桶的二阶项必然是非频繁的** 

e.g. hash(i,j, buckets\_len)=(i\*j) % buckets\_len





# 四验收流程



### 四 验收流程



- 检查1~4阶频繁项集和关联规则。
- 检查频繁项集和关联规则的数量。
- 提问了解编程思路和对Apriori算法和PCY算法的理解。