
电子科技大学计算机科学与工程学院

标准实验报告

(实验) 课程名称 数据结构与算法

电子科技大学

实验报告

学生姓名：陶浩轩 学号：2023080902011 指导教师：陈端兵

一、实验室名称：

学知三组团 10 栋 139

二、实验项目名称：

频繁模式挖掘

三、实验原理：

实验基于 Apriori 算法，利用先验性质，通过连接和剪枝操作，从给定的数据集中挖掘出频繁项集和关联规则。

四、实验目的：

- 理解频繁模式挖掘的概念和重要性。
- 掌握 Apriori 算法的基本原理和实现方法。
- 利用 Apriori 算法挖掘数据集中的频繁项集和关联规则。
- 分析频繁模式挖掘的结果，并进行解释和应用。

五、实验内容：

数据准备：随机生成指定范围的数据

参数设置：设置支持度阈值，确定挖掘频繁项集的最小支持度。

算法实现：编写 Apriori 算法的代码，实现频繁项集的挖掘。

结果分析：分析挖掘出的频繁项集和关联规则，并进行解释和应用。

六、实验器材（设备、元器件）：

平台: WSL Ubuntu 22.04 on Windows 10

七、实验步骤:

代码实现在 <https://github.com/WhaleFall-UESTC/UESTC-DataStructure/tree/master/proj1>

1. 定义 file, 随机生成数据以及将文件读取为 database, 并设置基本参数(随机变量以及最小支持度的范围)

```
01. #define OUTPUT_ALL 1
02. #define OUTPUT_MAX 2
03. #define OUTPUT_TOPN 4
04. #define OUTPUT_DEFAULT OUTPUT_ALL
05.
06. #define DEFAULT_FILE "src.txt"
07. #define DEFAULT_TOPN 5
08. #define DEFAULT_MIN_SUP 3
09.
10. #define RAND(1, u) (rand() % (1 - u + 1) + 1)
11. #define MAX_LEN 54
12. #define MIN_LEN 10
13. #define MAX_INDEX 128
14. #define MIN_INDEX 64
15. #define MAX_GROUP_NUM 64
16. #define MIN_GROUP_NUM 28
17.
18. #define BUF 512
19.
20. int read_file(database *db, const char *filename);
21. int random_file(const char *filename);
22.
```

2. 定义项集的存储方式: 类似于 bitmap, 定义 itemset 结构, 描述从某一地址开始一定长度的比特串, 其中从开始的第 i 位的比特, 1 表示编号 i 的元素存在于该项集, 0 表示不存在, 并定义 freqitem 与 freqlist 结构, 以链表的形式存储 itemset

```
01. #define NBITS 8 // itemset 块的位数
02. #define NBYTES (NBITS >> 3)
03. #define IDX2SIZE(idx) (((idx) >> 3) + 1)
04. typedef unsigned char uint8;
05. typedef uint8* itemset;
06. typedef uint8 itemblock;
07.
08. #define CHECKBIT(n, off) ((n) & (1 << off))
09. #define NO(i, j) (i * NBITS + j + 1)
10.
11. #define ALLOC 1
12. #define NOALLOC 0
13.
14. itemset init_itemset(int size);
15. int contain(itemset p, itemset c, int size);
16. itemset conjunct(itemset s1, itemset s2, int size, int k, int alloc);
17. void add_itemset(itemset items, int no);
18. void del_itemset(itemset items, int no);
19. itemset copy_itemset(itemset items, int size);
20. void print_itemset(itemset items, int size);
21.
```

```

01. struct freqitem {
02.     itemset items;
03.     int sup;
04.     struct freqitem* next;
05. };
06. typedef struct freqitem freqitem;
07.
08. struct freqlist {
09.     int len;        // length of List
10.     int size;       // size of itemset
11.     int k;
12.     freqitem list;
13. };
14. typedef struct freqlist freqlist;
15.
16. freqlist*      init_freqlist(freqitem* fl, int k, int size);
17. void          insert_freqlist(freqitem* fl, itemset items, int sup);
18. void          print_freqlist(freqitem* fl);
19. void          free_freqlist(freqitem* fl);
20. int          itemset_in_freqlist(freqitem* fl, itemset items);
21.

```

3. 定义 database 结构，此结构存储最初始的项集，是文件中读取的数据存储在内存中的结构

```

01. struct dbitem {
02.     itemset items;
03.     struct dbitem* next;
04. };
05. typedef struct dbitem dbitem;
06.
07. struct database {
08.     int len;
09.     int max_index;
10.     int size;
11.     dbitem list;
12. };
13. typedef struct database database;
14.
15. database*      init_db(database* db);
16. void          insert_db(database* db, itemset items);
17. int          count_minsup(database* db, itemset items);
18. void          print_db(database* db);
19. freqlist*      scan_db(database* db, int minsup);
20. void          free_db(database* db);
21.

```

4. 编写 debug.h 与 Makefile 方便调试与测试

5. 基础设施准备完毕，单元测试各模块后，开始在 main 函数实现 Apriori 算法

定义全局变量 db 以及与结果输出，文件读取相关的参数：

```

13. #define BUF 512
14. #define RING_NEXT(p) (((p) + 1) % RING_BUF)
15. #define RING_PRE(p) (((p) + RING_BUF - 1) % RING_BUF)
16.
17. #define LWC_SUP 1 // filter whose sup is less than min_sup
18. #define LWC_CHILD 2 // filter whose child is not frequent
19. #define LWC_REPEAT 4 // filter those repeat
20. #define LWC_ALL 7 // filter all
21.
22. #define SUP(sup) ((double) 1.0 * (sup) / db.len)
23.
24. static unsigned outset = OUTPUT_DEFAULT;
25. static char filename[BUF] = DEFAULT_FILE;
26. static int topn = DEFAULT_TOPN;
27. static bool rand_file = false;
28. static bool rand_only = false;
29. static int min_support = DEFAULT_MIN_SUP;
30.
31. static database db;

```

Apriori 算法的连接与剪枝操作。虽然从概念上这两个操作是解耦的，但是一些过滤需要二者的过程信息，出于效率考虑，我选择将它们合并在一起

```
74. freqlist*
75. link_with_cut(freqlist* fl, int flag)
76. {
77.     freqlist* ret = init_freqlist(NULL, fl->k + 1, fl->size);
78.     freqitem* pi = fl->list.next;
79.     bool enter_cycle = false;
80.     while (pi && pi->next) {
81.         enter_cycle = true;
82.         freqitem* pj = pi->next;
83.         while (pj) {
84.             itemset comb = conjunct(pi->items, pj->items, fl->size, fl->k + 1, ALLOC);
85.             if (comb == NULL) {
86.                 pj = pj->next;
87.                 continue;
88.             }
89.             bool pass_sup = true, pass_child = true;
90.             int sup = count_minsup(&db, comb);
91.             if (flag & LWC_SUP) {
92.                 if (sup < min_support) {
93.                     pass_sup = false;
94.                 }
95.             }
96.             if ((flag & LWC_CHILD) && pass_sup) {
97.                 int* nos = (int*) malloc(fl->size * sizeof(int));
98.                 int nop = 0;
99.                 int len = fl->size / NBYTES;
100.                 for (int i = 0; i < len; i++) {
101.                     for (int j = 0; j < NBYTES; j++)
102.                         if (CHECKBIT(comb[i], j)) {
103.                             nos[nop++] = NO(i, j);
104.                         }
105.                 }
106.                 for (int i = 0; i < nop; i++) {
107.                     itemset child = copy_itemset(comb, fl->size);
108.                     del_itemset(child, nos[i]);
109.                     if (itemset_in_freqlist(fl, child) < 0) {
110.                         free(child);
111.                         pass_child = false;
112.                         break;
113.                     }
114.                     free(child);
115.                 }
116.                 free(nos);
117.             }
118.             if (pass_sup && pass_child) {
119.                 insert_freqlist(ret, comb, sup);
120.             } else {
121.                 free(comb);
122.             }
123.             pj = pj->next;
124.         }
125.         pi = pi->next;
126.     }
127.
128.     if (!enter_cycle || ret->len == 0) {
129.         free(ret);
130.         return NULL;
131.     }
132.     return ret;
133. }
```

在 main 函数中，先解析参数，生成随机数据，读取数据并建立 database，并以此获取第一个 $k = 1$ 的频繁项集，随后不断调用 link_with_cut 筛选

```

135. int
136. main(int argc, char *argv[])
137. {
138.     init_db(&db);
139.     prase_args(argc, argv);
140.
141.     if (rand_file)
142.         random_file(filename);
143.     if (rand_only)
144.         return 0;
145.     read_file(&db, filename);
146.
147.     freqlist* fl[BUF];
148.     fl[0] = scan_db(db, min_support);
149.     print_freqlist(fl[0]);
150.     fl[1] = link_with_cut(fl[0], LWC_SUP);
151.     if (fl[1] == NULL) goto end;
152.     int j = 2;
153.     while ((fl[j] = link_with_cut(fl[j - 1], LWC_ALL)) != NULL)
154.         j++;
155.
156.     Log("Apriori Screening Finish");
157.

```

后续再处理输出结果。为了按照从小到大的顺序输出给定规则的频繁项集，我构建另一个变长的最小堆来存储 Apriori 每一轮的结果

```

01. #define HEAP_INIT 8
02.
03. struct heapitem {
04.     itemset items;
05.     int sup;
06. };
07. typedef struct heapitem heapitem;
08.
09. struct heap {
10.     heapitem *pq;
11.     int n;
12.     int capacity;
13. };
14. typedef struct heap heap;
15.
16. typedef unsigned char bool;
17. #define true 1
18. #define false 0
19.
20. heap*      init_heap();
21. void      insert_heap(heap* h, itemset items, int sup);
22. heapitem  pop_heap(heap* h);
23. void      add_freqlist(heap* h, freqlist *fl);
24. void      free_heap(heap* h);
25.

```

至此，初步实现了 Apriori 算法

八、实验数据及结果分析：

我生成了一个小的数据集来进行测试，设置 `min_support = 0.4`

```
2 3 4 6 7 9 10 11 12 15
1 6 7 9 10 11 12
1 4 5 8 9 10 12
1 3 4 6 8 10 11 12 13 14
1 6 8 9 11
2 3 7 8 10 11 12 13 15
3 5 7 8 9 11 13 14 15
1 3 4 5 7 13 14 15
5 8 11 14 15
2 3 4 5 7 11 14
```

这些数据可以表示一个图书馆内的部分顾客借阅过的书籍的类型

通过运行程序，我们可以得到一下的结果：

```
[main.c,158,main] Apriori Screening Finish
Maximal Frequent Pattern
sup      sets
0.40000  [ 3 7 11 ]
0.40000  [ 10 11 12 ]
0.40000  [ 3 7 15 ]

Most 5 Frequent Pattern
sup      sets
0.80000  [ 11 ]
0.60000  [ 3 ]
0.60000  [ 7 ]
0.60000  [ 8 ]
0.50000  [ 5 ]

All Frequent Pattern
sup      sets
0.80000  [ 11 ]
0.60000  [ 3 ]
0.60000  [ 7 ]
0.60000  [ 8 ]
0.50000  [ 5 ]
0.50000  [ 9 ]
0.50000  [ 3 7 ]
0.50000  [ 3 11 ]
0.50000  [ 10 ]
```

我们可以通过 Apriori 算法的到频繁项集，可以得知有哪些编号的关联性大

比如 3 类型代表犯罪心理学书籍，7 代表侦探小说，我们可以发现这两类书被同一人借阅的概率大，于是我们可以将这两类书的书架搬近一点

九、实验结论：

Apriori 算法是数据挖掘方面一个经典的算法。通过合并，剪枝，不断地筛选来找到频繁项集。我们可以通过这个方法来获取大量数据中的频繁模式，但是其平方时间复杂度使得其效率有些堪忧（一些较大但有没那么大的数据集甚至跑到了十几秒）。在使用的过程中，我们

应当多采取一下优化的措施，改良算法

十、总结及心得体会：

1. 本次我使用了 bitmap 的形式存储频繁项集，不仅节省了空间（虽然 Apriori 更关心时间复杂度），而且在判断是否为子集时只需要将二者 and，合并时只需要将二者 or，操作上也更加快捷
2. 我尝试设计了一套层级结构：freqitem 与 freqlist 建立在 itemset 之上；database 建立在 freqlist 之上，这套层级结构方便了程序的构建

十一、对本实验过程及方法、手段的改进建议：

1. 我们需要在 database 中查找频繁项集出现的次数，采用的是顺序查找的方式。优化的方案是：我们可以利用哈希表来减小查找次数。根据 itemset 的长度，我们可以按照偏移量均匀地选取 4 个比特，并用这 4 个比特所组成的无符号数的值作为 hashcode
2. 如果 itemset 的长度较小（小于 128），那么这个 itemset 是对 C 而言一个天然的无符号数表示。我们可以在哈希桶内采用二分查找（或者构建二叉搜索树）的方式来减小搜索次数。此外，还可以根据 itemset 的长度与数量的大小决定使用基数排序还是快速排序
3. 对应的，在 database 中可以采用变长数组的方式而非链表（即初始长度固定，若所要的容量超过了这个长度，则再开辟一个大小为原来的两倍的空间，将数据复制到这个空间，释放原来的空间），可以加快查找速度
4. 或者说构建一棵树，以一个项集的第一个元素作为根，若是其他项集有与其共同的部分，但是后续的项目有所不同，则在这棵树上“分支”，并且树要求子节点必须按照某一顺序大于根节点。这样在计数的时候可以快速地获取支持度
5. 最开始写代码的时候可以参照网上的测试集的格式来编写
6. 再分析优化程序时，可以使用 `ftrace` 跟踪函数调用与开销，针对性地优化

报告评分：

指导教师签字：