



西安电子科技大学
XIDIAN UNIVERSITY

数据挖掘第二次上机实验报告

基于 Apriori 算法的国会投票关联规则挖掘

姓名：刘臻劼

学号：19030700016

班级：1903071

导师：郭杏莉

目录

1	实验内容	2
1.1	实验目标	2
1.2	数据集介绍	2
2	实验分析	3
2.1	题目分析与数据预处理	3
2.2	先验原理与 Apriori 算法概述	3
3	具体实现	4
3.1	数据预处理	4
3.2	挖掘 1-频繁集与 2-频繁集	6
3.3	挖掘所有频繁集	8
3.4	生成第一类关联规则	8
3.5	生成第二类关联规则	9
4	实验结果	10
4.1	第一类关联规则	10
4.2	第二类关联规则	10
5	结语与心得	11

1 实验内容

1.1 实验目标

本次实验目标为：使用 Apriori 算法，支持度设为 30%，置信度为 90%，挖掘高置信度的规则。

1.2 数据集介绍

Congressional Voting Records 数据集包含 435 条议员对于 16 项提案的投票记录（其中含 267 条民主党投票与 168 条共和党投票）。

每条记录包含 17 项，其中第一项为党派：republican 或 democrat；其余 16 项是对于各个议题的投票情况，取值为 y、n、?，分别表示支持、反对与弃权。

以下为各提案弃权票情况。

表 1: 各项提案的弃权票

提案	弃权票数
1.handicapped-infants	0
2.water-project-cost-sharing	12
3.adoption-of-the-budget-resolution	48
4.physician-fee-freeze	11
5.el-salvador-aid	11
6.religious-groups-in-schools	15
7.anti-satellite-test-ban	11
8.aid-to-nicaraguan-contras	14
9.mx-missile	15
10.immigration	22
11.synfuels-corporation-cutback	7
12.education-spending	21
13.superfund-right-to-sue	31
14.crime	25
15.duty-free-exports	17
16.export-administration-act-south-africa	28

2 实验分析

2.1 题目分析与数据预处理

(1) 数据预处理

我们可将党派与各提案投票情况看做一个项，各项列举如下：

表 2: 各项提案的弃权票

项	项
I_1 .Republican	I_2 .Democrat
I_3 .handicapped-infants=Y	I_4 .handicapped-infants=N
I_5 .water-project-cost-sharing=Y	I_6 .water-project-cost-sharing=N
I_7 .adoption-of-the-budget-resolution=Y	I_8 .adoption-of-the-budget-resolution=N
I_9 .physician-fee-freeze=Y	I_{10} .physician-fee-freeze=N
I_{11} .el-salvador-aid=Y	I_{12} .el-salvador-aid=N
I_{13} .religious-groups-in-schools=Y	I_{14} .religious-groups-in-schools=N
I_{15} .anti-satellite-test-ban=Y	I_{16} .anti-satellite-test-ban=N
I_{17} .aid-to-nicaraguan-contras=Y	I_{18} .aid-to-nicaraguan-contras=N
I_{19} .mx-missile=Y	I_{20} .mx-missile=N
I_{21} .immigration=Y	I_{22} .immigration=N
I_{23} .synfuels-corporation-cutback=Y	I_{24} .synfuels-corporation-cutback=N
I_{25} .education-spending=Y	I_{26} .education-spending=N
I_{27} .superfund-right-to-sue=Y	I_{28} .superfund-right-to-sue=N
I_{29} .crime=Y	I_{30} .crime=N
I_{31} .duty-free-exports=Y	I_{32} .duty-free-exports=N
I_{33} .export-administration-act-south-africa=Y	I_{34} .export-administration-act-south-africa=N

则有全项集 $I = \{I_1, \dots, I_{34}\}$ 。则可将每项记录转换为事务，有事务集 $T = \{T_1, \dots, T_{435}\}$ 。其中： $I(T_i) \subseteq I, T_i \in T$ 。

(2) 挖掘目标鉴于强关联规则数量过多，此处考虑挖掘可解释性较好的两类关联规则：

1. 根据投票情况挖掘党派信息：形如 $X \Rightarrow Y$ 的关联规则，其中 $X \subseteq I - \{I_1, I_2\}$, $Y = I_1$ 或 $Y = I_2$ ；

2. 某党派内部的投票规则：形如 $X \Rightarrow Y$ 的关联规则，其中 $X, Y \subseteq I - \{I_1, I_2\}$, X 包含且仅包含 I_1 、 I_2 中的一个；

2.2 先验原理与 Apriori 算法概述

对于频繁项集与其子集，我们有如下的先验原理：

一个项集是频繁的，则其所有子集是频繁的；一个项集是非频繁的，则其所有超集是非频繁的。

该原理启示我们可以在项集格上采用基于支持度的剪枝策略以降低候选项集的数量。由此我们有如下算法。记 C_k 为 k -候选集， F_k 为 k -频繁项集。

其中通过 $(k-1)$ -频繁项集产生 k -候选项集的算法 $C_k_Generator$ 采用 $F_{k-1} \times F_{k-1}$ 策略。

Algorithm 1 基于支持度剪枝的频繁项集产生算法 (Apriori)

Output: 所有频繁项集的集合 F 。

```
1:  $F_1 = \{i | i \in I \wedge \sigma(i)/|I| \geq SUP_{min}\}, k = 2$ 
2: while  $F_k \neq \Phi$  do
3:    $C_k = C_k\_Generator(F_{k-1})$ 
4:   for  $t \in T$  do
5:     for  $c \in C_k$  do
6:       if  $issubset(c, t)$  then
7:          $\sigma(c) = \sigma(c) + 1$ 
8:       end if
9:     end for
10:  end for
11:   $F_k = \{c | c \in C_k \wedge \sigma(c)/|I| \geq SUP_{min}\}$ 
12:   $k = k + 1$ 
13: end while
14:  $F = \bigcup F_k$ 
15: return  $F$ ;
```

Algorithm 2 $C_k_Generator$ 算法

Output: (k-1)-频繁项集 F_{k-1} 。

```
1:  $F_k = \{\}$ 
2: for  $c_1 \in F_{k-1}$  do
3:   for  $c_2 \in F_{k-1}$  do
4:     if  $c_1^{(i)} = c_2^{(i)}, i = 1, 2, \dots, k-2 \wedge c_1^{k-1} \neq c_2^{k-2}$  then
5:        $F_k = F_k + \{c_1^{(i)}, c_2^{k-2}\}, i = 1, 2, \dots, k-1$ 
6:     end if
7:   end for
8: end for return  $F_k$ 
```

3 具体实现

3.1 数据预处理

首先读取并加载数据：

```
1 SUP_min = 0.3 # 最小支持度
2 CONF_min = 0.9 # 最小置信度
3 # 加载数据，构建事务集
4 def loadData():
5     f = open('data/house-votes-84.data')
6     data = []
7     for txt_line in f.read().split('\n'):
8         data_processed = []
9         if txt_line == '':
10             pass
11         else:
12             tmp = txt_line.split(',')
```

```

13         data_processed.append(1) if tmp[0] == 'republican' else
           data_processed.append(2)
14     for i in range(len(tmp)):
15         if tmp[i] == 'y':
16             data_processed.append(2 * i + 1)
17         elif tmp[i] == 'n':
18             data_processed.append(2 * i + 2)
19         else:
20             pass
21     data.append(data_processed)
22 f.close()
23 return data
24 raw_data = loadData()

```

根据任务分割数据集:

```

1 def split_data(raw_data, party):
2     if party == '':
3         return raw_data
4     elif party == 'republic':
5         data_republic = []
6         for data in raw_data:
7             if len(data) == 1:
8                 continue
9             if data[0] == 1:
10                 data_republic.append(data[1:])
11         return data_republic
12     elif party == 'democrat':
13         data_democrat = []
14         for data in raw_data:
15             if len(data) == 1:
16                 continue
17             if data[0] == 2:
18                 data_democrat.append(data[1:])
19         return data_democrat
20 #data_republic = split_data(raw_data, party='republic')

```

```

21 #data_democrat = split_data(raw_data, party='democrat')
22 data = split_data(raw_data, party='')

```

3.2 挖掘 1-频繁集与 2-频繁集

```

1 def compare(item1, item2):
2     #判断两个 (k-1) 候选项集前 (k-2) 项是否相等
3     tmp = []
4     equal = True
5     if len(item1) == 1:
6         return equal, (item1 + item2)
7     for i in range(len(item1) - 1):
8         tmp.append(item1[i])
9         if item1[i] != item2[i]:
10             equal = False
11             break
12     if equal:
13         minitem = min(item1[-1], item2[-1])
14         maxitem = max(item2[-1], item1[-1])
15         tmp.append(minitem)
16         tmp.append(maxitem)
17     return equal, tuple(tmp)
18
19 def issubset(item1, item2):
20     # 判断 item1 是否为 item2 的子集
21     return set(item1).issubset(set(item2))
22
23 #挖掘 1-频繁集
24 def F1_Mining(data):
25     C1, F1 = {}, {}
26     count = len(data)
27     for data_line in data:
28         # 对于数据集中的每一行投票数据
29         for i in range(len(data_line)):
30             # 对于每一行数据中的下标 (对应某个议题)

```

```

31         key = tuple([data_line[i]])
32         if key in C1.keys():
33             # 以键值对的形式进行存储和计数
34             C1[key] += 1
35         else:
36             C1[key] = 1
37     for item in C1:
38         if C1[item] / count >= SUP_min:
39             F1[item] = C1[item]
40     return F1
41
42 # 挖掘 2-频繁集
43 def F2_Mining(data, F1):
44     C2, F2 = {}, {}
45     F1_keys = sorted(list(F1.keys()))
46     for i in range(len(F1_keys)):
47         for j in range(i + 1, len(F1_keys)):
48             equal, tmp = compare(F1_keys[i], F1_keys[j])
49             if equal:
50                 C2[tmp] = 0
51             else:
52                 continue
53     for data_line in data:
54         for item in C2:
55             if issubset(item, data_line):
56                 C2[item] += 1
57     for item in C2:
58         if C2[item] / len(data) >= SUP_min:
59             F2[item] = C2[item]
60     return F2
61
62 F1 = F1_Mining(data)
63 F2 = F2_Mining(data, F1)

```


3.3 挖掘所有频繁集

```
1 # 挖掘3-频繁项集乃至n-频繁项集，并合并所有频繁项集
2 def F_Mining(data, F1, F2):
3     Fk_1 = F2
4     F = [F1] # 存储的是所有频繁项集
5     while Fk_1 != {}:
6         F.append(Fk_1)
7         Ck, Fk = {}, {}
8         Fk_1_keys = sorted(list(Fk_1.keys()))
9         for i in range(len(Fk_1_keys)):
10             for j in range(i + 1, len(Fk_1_keys)):
11                 equal, tmp = compare(Fk_1_keys[i], Fk_1_keys[j])
12                 if equal:
13                     Ck[tmp] = 0
14                 else:
15                     continue
16         for data_line in data:
17             for item in Ck:
18                 if issubset(item, data_line):
19                     Ck[item] += 1
20         for item in Ck:
21             if Ck[item] / len(data) >= SUP_min:
22                 Fk[item] = Ck[item]
23         Fk_1 = Fk
24     return F
25
26 F = F_Mining(data, F1, F2)
```

3.4 生成第一类关联规则

```
1 def generate_1_rules(F):
2     rules = {}
3     for i in range(2, len(F)):
4         for item in F[i]:
5             if 1 in item:
```

```

6         antecedent = tuple(sorted(set(item) - {1})) # 关联规则
           前件
7         confidence = F[i][item] / F[i - 1][antecedent]
8         if 1 > confidence >= CONF_min:
9             rules[(antecedent, 1)] = confidence
10        elif 2 in item:
11            antecedent = tuple(sorted(set(item) - {2})) # 关联规则
               前件
12            confidence = F[i][item] / F[i - 1][antecedent]
13            if 1 > confidence >= CONF_min:
14                rules[(antecedent, 2)] = confidence
15        else:
16            continue
17    return rules
18
19 rules = generate_1_rules(F).items()
20 rules = sorted(rules, key=lambda x: x[1], reverse=True)

```

3.5 生成第二类关联规则

```

1 def get_subset(freq):
2     # 获得某一元组的所有子集
3     if len(freq) == 0:
4         return [[]]
5     all_subsets = get_subset(freq[1:]) + [[freq[0]] + r for r in
        get_subset(freq[1:])]
6     all_subsets = sorted(all_subsets, key=lambda x: (len(x), x))
7     return all_subsets
8
9 def subset_tuple(all_subsets, freq):
10    # 整理子集
11    all_subsets_tuple = []
12    for subset in all_subsets:
13        if len(subset) == 0 or len(subset) == len(freq):
14            continue

```

```

15         else:
16             subset = set(subset)
17             complement = set(list(freq)) - subset
18             subset = tuple(sorted(subset, key=lambda x: x))
19             complement = tuple(sorted(complement, key=lambda x: x))
20             all_subsets_tuple.append(tuple([subset, complement]))
21     return all_subsets_tuple
22
23 def generate_2_rules(F, party=party):
24     # 生成第二类关联规则
25     rules = {}
26     for i in range(1, len(F)):
27         for item in F[i]:
28             for subset in subset_tuple(get_subset(item), item):
29                 rule_confidence = F[len(item) - 1][item] / F[len(subset
30                                     [0]) - 1][subset[0]]
31                 if rule_confidence >= CONF_min:
32                     rules[subset] = rule_confidence
33     return rules
34
35 rules = sorted(generate_2_rules(F).items(), key=lambda x: x[1], reverse
36               =True)

```

4 实验结果

4.1 第一类关联规则

第一类关联规则即根据投票情况挖掘党派信息，共挖掘出 202 条关联规则。高置信度的规则如下（部分）：

```

physician-fee-freeze=N, aid-to-nicaraguan-contras=Y, -> Democrat, 置信度:0.99526
physician-fee-freeze=N, education-spending=N, -> Democrat, 置信度:0.99505
physician-fee-freeze=N, el-salvador-aid=N, -> Democrat, 置信度:0.99487
physician-fee-freeze=N, el-salvador-aid=N, aid-to-nicaraguan-contras=Y, -> Democrat, 置信度:0.99479
physician-fee-freeze=N, anti-satellite-test-ban=Y, aid-to-nicaraguan-contras=Y, -> Democrat, 置信度:0.99459
physician-fee-freeze=N, aid-to-nicaraguan-contras=Y, education-spending=N, -> Democrat, 置信度:0.99448
physician-fee-freeze=N, mx-missile=Y, -> Democrat, 置信度:0.99444
physician-fee-freeze=N, el-salvador-aid=N, anti-satellite-test-ban=Y, -> Democrat, 置信度:0.99438
physician-fee-freeze=N, el-salvador-aid=N, anti-satellite-test-ban=Y, aid-to-nicaraguan-contras=Y, -> Democrat, 置信度:0.99429
physician-fee-freeze=N, aid-to-nicaraguan-contras=Y, mx-missile=Y, -> Democrat, 置信度:0.99425

```

4.2 第二类关联规则

第一类关联规则即党派内部投票规律。共和党内部分高置信度的投票规律如下：

```
关联规则:water-project-cost-sharing=Y, -> physician-fee-freeze=Y, 置信度:1.00000
关联规则:water-project-cost-sharing=N, -> crime=Y, 置信度:1.00000
关联规则:handicapped-infants=N, water-project-cost-sharing=Y, -> physician-fee-freeze=Y, 置信度:1.00000
关联规则:handicapped-infants=N, water-project-cost-sharing=N, -> crime=Y, 置信度:1.00000
关联规则:handicapped-infants=N, anti-satellite-test-ban=N, -> physician-fee-freeze=Y, 置信度:1.00000
关联规则:handicapped-infants=N, immigration=Y, -> physician-fee-freeze=Y, 置信度:1.00000
关联规则:handicapped-infants=N, education-spending=Y, -> physician-fee-freeze=Y, 置信度:1.00000
关联规则:handicapped-infants=N, superfund-right-to-sue=Y, -> physician-fee-freeze=Y, 置信度:1.00000
关联规则:handicapped-infants=N, duty-free-exports=N, -> physician-fee-freeze=Y, 置信度:1.00000
关联规则:water-project-cost-sharing=Y, adoption-of-the-budget-resolution=N, -> physician-fee-freeze=Y, 置信度:1.00000
```

民主党内部分高置信度的投票规律如下:

```
关联规则:handicapped-infants=Y, synfuels-corporation-cutback=N, -> physician-fee-freeze=N, 置信度:1.00000
关联规则:water-project-cost-sharing=N, el-salvador-aid=N, -> physician-fee-freeze=N, 置信度:1.00000
关联规则:water-project-cost-sharing=N, anti-satellite-test-ban=Y, -> physician-fee-freeze=N, 置信度:1.00000
关联规则:water-project-cost-sharing=N, aid-to-nicaraguan-contras=Y, -> physician-fee-freeze=N, 置信度:1.00000
关联规则:water-project-cost-sharing=N, superfund-right-to-sue=N, -> physician-fee-freeze=N, 置信度:1.00000
关联规则:water-project-cost-sharing=N, export-administration-act-south-africa=Y, -> physician-fee-freeze=N, 置信度:1.00000
关联规则:el-salvador-aid=N, synfuels-corporation-cutback=N, -> physician-fee-freeze=N, 置信度:1.00000
关联规则:religious-groups-in-schools=N, synfuels-corporation-cutback=N, -> physician-fee-freeze=N, 置信度:1.00000
关联规则:anti-satellite-test-ban=Y, synfuels-corporation-cutback=N, -> physician-fee-freeze=N, 置信度:1.00000
关联规则:aid-to-nicaraguan-contras=Y, synfuels-corporation-cutback=N, -> physician-fee-freeze=N, 置信度:1.00000
```

注意到前 10 的关联规则置信度均为 1, 这是由于数据集较小, 党派内在某几项议题上的投票情况完全相同。

5 结语与心得

本实验中统计候选项集的支持度时还是暴力遍历方法, 没有使用教材上的哈希树等方法, 算是偷了点小懒, 今后应当更加努力, 尽可能优化算法。