

# 基于类关联规则的分类器的实现

刘砺志                  陈大松                  肖璐菁

22920142203873    22920142203785    22920142203933

{22920142203873, 22920142203785, 22920142203933}@stu.xmu.edu.cn

**摘 要：**分类问题就是确定对象属于哪一个预定义的目标类的过程。传统的关联分析通过挖掘隐藏在大型数据集背后的满足最小支持度和最小置信度约束的关联规则，来描述数据集中存在的有意义的联系，但是所发掘的联系都是事先无法预知的。通过将分类任务和关联分析相结合，我们利用一个关联规则的特殊子集，称为类关联规则，来训练一个分类器，对记录进行分类。实验结果表明，使用这种技术训练出的分类器，其准确率一般来说都是很高的。

**关键字：**分类；类关联规则；关联分析

## An Implementation of Classifier Based on Class Association Rules

Lizhi Liu                  Dasong Chen                  Lujing Xiao

22920142203873    22920142203785    22920142203933

{22920142203873, 22920142203785, 22920142203933}@stu.xmu.edu.cn

**Abstract:** The classification problem is the process of determining which pre-defined target class the object belongs to. The traditional association analysis describes the meaningful links in the data set by mining the association rules that satisfy the minimum support and the minimum confidence constraint behind the large data set, but those links are unpredictable in advance. By combining classification task with association analysis, we use a special subset of association rules, called class association rules, to train a classifier to classify records. Experimental results show that the use of this technology training out of the classifier, the accuracy rate is very high in general.

**Keywords:** Classification; Class Association Rules; Association Analysis;

# 1 引言

传统的关联规则挖掘意在从海量数据中挖掘出有意义的规则，以描述数据之间存在的相互联系。而分类问题则是利用分类器，判断一个样本究竟应该归为哪一个预定义的目标类。对于前者，所能发现的规则是无法事先知晓的；对于后者，目标类已经事先划定并且最终的判定结果只能有一个。因而合理的整合两种技术，以构造一个更为强大的分类器，成为了我们这篇文章所要实现的东西。本文中通过修改 Apriori 算法 [3]，寻找一组特殊的频繁项集，以生成一组特殊的关联规则，称为类关联规则。我们可以使用这些类关联规则，去训练一个基于规则的分类器，实现对样本的分类。

由于现实中存在的数据集完整程度高低有别，各属性值类型也千差万别，在挖掘类关联规则之前，我们必须要对数据进行预处理。预处理过程主要包括两个部分。一个是对缺失值进行处理。对于缺失值较多的特征，我们直接舍弃；对于缺失值较少的特征，我们可以把缺失本身直接作为一个特征，或是使用均值、上下文数据、众数法进行填充，也可以使用插值、随机森林算法 [4] 拟合等方式填补。第二个问题是对连续型属性进行离散化的问题。文献 [5] 和 [6] 提出了基于信息熵的连续型属性值离散化方法，文献 [7] 中对这项工作进行了很好的概述。我们将在第 2 部分第 1 节详细介绍我们所采用的方法。

本文的重点和难点在于如何生成类关联规则并利用这些规则训练出性能优异的分类器。文献 [1] 提出的 CBA 算法很好的解决了这个问题。这个算法包含两个部分。第一部分是规则生成，即 CBA-CG 算法。这个算法在 Apriori 算法的基础上进行修改，以生成一组类关联规则而非传统的关联规则。这一部分实现的关键在于如何筛选出候选项集，以及如何利用悲观误差进行规则剪枝。这一部分我们将在第 2 部分第 2 节进行讨论。算法的第二部分就是训练分类器。刘兵等人首先提出了一个简单的 CBA-CB M1 算法，利用贪心策略，按序逐步抽取规则，最后构建一个有序规则集来作为最终的分类器。这一部分我们将在第 2 部分第 3 节进行讨论。但是这个算法对于大样本情形并不适用，因此又提出了改进版的 CBA-CB M2 算法。这个算法包括 3 个阶段。第 1 阶段对类关联规则进行分类，缩小最终用于生成分类器的规则数。第 2 阶段对未决规则进行裁决，确定最后用于分类器生成的规则。第 3 阶段则在筛选出的那一部分规则中提取用于分类器的那些规则。这一部分的详细讨论将在第 2 部分第 4 节进行。

完成了分类器的生成，就需要对其性能进行检验。我们选取了 UCI 机器学习知识库 [2] 中的 30 个数据集，进行了 10 折交叉检验 [8]。通过检验结果分析，我们所实现的分类器的判定精度极高。这一部分的分析讨论详见本文第 3 部分。

## 2 方法

### 2.1 数据预处理

从实际中采集到的数据往往存在着格式不统一、属性值缺失、噪声点等一系列问题。尽管我们实现时所使用的数据都是从 UCI 机器学习资源库中下载的，这些数据已经尽管了相关人员的一些预处理，但是仍然存在两个主要问题：a) 存在缺失值和 b) 连续型属性值离散化。下面我们将就这两个方面进行介绍。

#### 2.1.1 缺失值处理

一般来说，对于缺失值，我们要分两种情况去考虑。对于缺失值较多的属性，我们应当毫无保留的删去，因为这样属性的存在将会对最后分类器的性能带来极大影响。对于缺失值较少的属性，则应当考虑保留。我们可以采取使用均值、众数、上下文数据（即紧邻该样本的上一条或下一条不含缺失值的记录）进行填充这样简单的策略，也可以使用插值或是使用随机森林方法进行拟合，完成数据的填补。

在本文中，由于我们的工作重心是在 CBA 算法的实现而非数据的预处理，因此我们选取了简单的策略：对于缺失值较少的特征，利用缺失值所属特征的众数进行填补；对于缺失值较多的特征，则丢弃这一列。这里选取的缺失率阈值  $\varepsilon = 0.5$ 。即

$$m = \begin{cases} \text{mode}(\{x : x \in \text{col}[m] \wedge x \neq ?\}) & \text{ratio} < \varepsilon \\ \text{discard col}[m] & \text{otherwise} \end{cases}$$

上式中， $m$  表示缺失值； $\text{mode}(\cdot)$  操作表示取  $\cdot$  的众数（若存在多个众数，则随机选取一个）； $\text{col}[m]$  表示  $m$  所属的列； $?$  指代缺失值； $\text{ratio}$  表示缺失率，其计算方式为  $\text{ratio} = \#m/N$ ，其中  $\#m$  表示缺失值的个数， $N$  表示这一列涵盖缺失值下的总样本数。

### 2.1.2 连续型属性值的离散化

所谓连续性属性值，就是该属性的取值范围并非有限集。与之相对则是离散型的特征，其取值只可能落在一个有限集合中。在实际的数据集中，连续型属性值非常常见，比如商品的价格、某地的气温等等。但是在经典的决策树模型中，我们只能处理离散型的特征。因此我们需要对连续型属性进行离散化。例如对于商品的价格，在我们初始数据中，可能是具体的数字，但是经过我们的处理后，它就变成了便宜、中等、昂贵三个层次，用形式化的语言描述就是完成这样的映射：

$$\mathbb{R}^+ \rightarrow \{\text{cheap, medium, expensive}\}$$

文献 [7] 总结了文献 [5][6] 提出的递归的最小熵划分算法。它通过对该属性的取值范围不断进行划分，分隔出几个区间，每个区间对应于一个离散值，由此完成离散化。这个算法主要思想有两点：**a)** 使用熵去衡量一次划分的优劣，即这次划分会将数据集的信息量最大化，以及 **b)** 划分不能无限制进行下去，以避免过拟合的出现。

给定样本集  $S$ ，特征  $A$  以及划分边界  $T$ ，由边界  $T$  引起的划分下的类信息熵定义为

$$E(A, T; S) = \frac{|S_1|}{|S|} \text{Ent}(S_1) + \frac{|S_2|}{|S|} \text{Ent}(S_2) \quad (1)$$

对于给定的特征  $A$ ，在所有划分中能够使得上述熵值最小化的那个划分边界  $T_{\min}$ ，将被挑选出来作为这一轮划分的二元划分边界。事实上，这个定义就是将划分后产生的两个子集的信息熵进行加权平均，即考虑划分产生的两个子集的综合效果。在本文中，我们采用的信息熵定义为

$$\text{Ent}(S) = - \sum_{i=1}^n P(x_i) \log_2 P(x_i)$$

其中， $n$  表示样本集  $S$  中的类的种数， $P(x_i)$  表示样本集中第  $i$  类出现的概率，在实际计算时使用频率以近似替代。

但是作为递归程序，必须还要设计递归出口，否则划分就将无休止的持续下去。为此，Fayyad 和 Irani 提出了最小描述长度准则，来决定何时停止递归划分。这个准则告诉我们，对于集合  $S$  的递归划分停止，当且仅当

$$\text{Gain}(A, T; S) < \frac{\log_2(N-1)}{N} + \frac{\Delta(A, T; S)}{N} \quad (2)$$

成立。其中  $N$  是集合  $S$  的元素个数；信息增益  $\text{Gain}(A, T; S)$  定义为

$$\text{Gain}(A, T; S) = \text{Ent}(S) - E(A, T; S)$$

即划分前后的信息量的增量；余量

$$\Delta(A, T; S) = \log_2(3^k - 2) - [k \cdot \text{Ent}(S) - k_1 \cdot \text{Ent}(S_1) - k_2 \cdot \text{Ent}(S_2)]$$

其中  $k_i$  表示集合  $S_i$  中的样本个数,  $i = 1, 2$ 。不等式 (2) 事实上给出了信息增益的下界, 若一次划分的信息增益不能超过这个下界, 则不再继续往下进行划分。格式化的算法描述如算法 1 所示。

下面我们给出一个简单的例子 [9], 手工执行一下整个算法。表 1 给出了我们的样本集, 其中包括 5 个样本, 每个样本包含 2 个分量, 一个分量是连续型特征, 另一个是二元分类结果, 仅包括 Yes 和 No 两类。

表 1 用于解释递归的最小熵划分算法的简单数据集

Continuous Feature	Output Class
1.0	Yes
1.0	Yes
2.0	No
3.0	Yes
3.0	No

初始时, 信息熵  $\text{Ent}(S) = -\frac{3}{5} \times \log_2 \frac{3}{5} - \frac{2}{5} \times \log_2 \frac{2}{5} = 0.97$ 。下面开始选取二元分隔点。若选择 1.0 作为分界点, 则对于严格小于 1.0 的这一部分, 事实上为空, 而大于等于 1.0 的另一部分, 实际上就是原数据集, 因此并没有产生增益, 显然不满足式 (2) 的条件。类似地, 可以分析选择 3.0 作为分界点的情形。下面着重讨论以 2.0 作为分界点的情形。分隔完成后, 数据集被划分成两部分, 一部分为小于 2.0 的那一部分, 如表 4 所示, 另一部分为大于等于 2.0 的那一部分, 如表 3 所示。

表 2 对原始数据集以 2.0 作为分界点小于 2.0 的分块数据

Continuous Feature	Output Class
1.0	Yes
1.0	Yes

表 3 对原始数据集以 2.0 作为分界点大于等于 2.0 的分块数据

Continuous Feature	Output Class
2.0	No
3.0	Yes
3.0	No

由表 4 数据, 不难算出这一部分  $S_1$  的信息熵  $\text{Ent}(S_1) = -1 \times \log_2 1 = 0$ 。而另一部分  $S_2$  的信息熵  $\text{Ent}(S_2) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.92$ 。所以此处划分后的类信息熵  $E = \frac{2}{5} \text{Ent}(S_1) + \frac{3}{5} \text{Ent}(S_2) = 0.55$ 。因此此次划分产生的信息增益  $\text{Gain} = \text{Ent}(S) - E = 0.42$ 。根据余量公式, 可以求出  $\Delta = \log_2(3^5 - 2) - [5 \times \text{Ent}(S) - 2 \times \text{Ent}(S_1) - 3 \times \text{Ent}(S_2)] = 5.82$ 。由此, 可以计算出信息增益下界  $\inf \text{Gain} = \frac{\log_2(5-1)}{5} + \frac{\Delta}{5} = 1.56$ 。因为  $\text{Gain} = 0.42 < 1.56$ , 所以不进行划分。于是 Split 返回空集, Partition 也不再继续递归, 那么  $T$  为空, 即整个数据集的连续属性值全部标定成一类。

从这个简单的例子中我们发现了这个方法的一个问题, 如果这一列的数据信息量不足以将其划分若干列, 或者说划分的粒度过粗, 则会导致之后的规则生成和分类器训练产生偏差。因此, 我们在实际应用中,

---

**算法 1** 递归的最小熵划分算法

---

```
1:  $T = \emptyset$ 
2:
3: function Split( $S$ )
4:   sort( $S$ )
5:    $W = \emptyset$ 
6:   for each candidate  $c$  in  $S$  do
7:     calculate  $g(c) = \text{Gain}(A, c; S)$ 
8:     if (2) is not satisfied then
9:        $W = W \cup c$ 
10:    end if
11:  end for
12:  if  $W = \emptyset$  then
13:    return  $\emptyset$ 
14:  else
15:    return  $\arg \max_{c \in W} g(c)$ 
16:  end if
17: end function
18:
19: function Partition( $S$ )
20:    $t = \text{Split}(S)$ 
21:   if  $t = \emptyset$  then
22:     return
23:   else
24:      $T = T \cup t$ 
25:      $S_1 = \{x : x \in S \wedge x < t\}$ 
26:      $S_2 = \{x : x \in S \wedge x \geq t\}$ 
27:     Partition( $S_1$ )
28:     Partition( $S_2$ )
29:   end if
30: end function
31:
32: function main
33:   Partition( $S$ )
34:   sort( $T$ )
35:   return  $T$ 
36: end function
```

---

若最小熵划分算法返回的划分边界集合为空，则从简计议，简单的将这一属性的取值区间等分成 3 份，对于样本  $x \in [\min(S), \min(S) + \frac{\max(S) - \min(S)}{3}]$ ，将其标为第一块；样本  $x \in (\min(S) + \frac{\max(S) - \min(S)}{3}, \min(S) + 2 \times \frac{\max(S) - \min(S)}{3}]$ ，将其标为第二块；最后把  $x \in (\min(S) + 2 \times \frac{\max(S) - \min(S)}{3}, \max(S)]$  标为第三块。这样就避免了划分粒度过粗带来的影响。

## 2.2 规则生成的 CBA-CG 算法实现

### 2.2.1 规则项集和类关联规则

CBA-RG 算法的核心就是去寻找所有满足最小支持度约束的规则项集 (ruleitem)。因而在具体讨论算法之前，有必要对规则项集的相关定义进行一些讨论。

一个规则项集就是一个形如

$$\langle \text{condset}, y \rangle$$

的偶对，其中， $\text{condset}$  是项集的集合， $y \in Y$  是所有类标号中的一个类标号，它表示类关联规则

$$\text{condset} \rightarrow t$$

举例来说，表 4 给出了一个简单的数据集 [10]，其中包括 A、B 两个属性，以及 C 这一两类划分结果，每个样本都被划分入 y 或 n。比如， $\langle \{(A, e), (B, p)\}, (C, y) \rangle$  就是一个规则项集，它表示了类关联规则  $\{(A, e), (B, p)\} \rightarrow y$ 。从这里不难看出，事实上，类关联规则就是一个简单的 IF...THEN... 规则，前件  $\text{condset}$  就是一个合取范式，它反映了属性集的一个子集的取值情况，而后件  $y$  则表示如果前件成立时该样本应归为哪一类。就前述的例子而言，就可以写为  $(A = e) \wedge (B = p) \rightarrow y$ ，即当  $A = e$  并且  $B = p$  时，样本应该归为  $y$  类。

表 4 用于说明 CBA 算法的示例数据集

A	B	C
e	p	y
e	p	y
e	q	y
g	q	y
g	q	y
g	q	n
g	w	n
g	w	n
e	p	n
f	q	n

定义  $\text{condset}$  的支持度计数  $\text{condsupCount}$  为数据集  $D$  上包含有  $\text{condset}$  的样本个数，其实就是统计数据集中对应于  $\text{condset}$  的那些属性值的取值与  $\text{condset}$  的取值一致的那些样本的个数。就前述的规则项集而言，有 3 个样本的 A 和 B 属性取值满足  $\text{condset}$  的取值（即  $A = e, B = p$ ），因此  $\text{condsupCount} = 3$ 。

更进一步，定义规则项集的支持度计数  $\text{rulesupCount}$  为数据集  $D$  上包含有  $\text{condset}$  并且最后的类标号也一致的样本个数，事实上就是统计数据集中不仅满足  $\text{condset}$  的取值且最后分类一致的样本个数。还是看前面的例子，前面已经知道  $\text{condsupCount} = 3$ ，而其中最后归为  $y$  类的样本个数只有 2 个，因而  $\text{rulesupCount} = 2$ 。

根据前面的两个支持度定义，我们便可以定义类关联规则  $condset \rightarrow t$  的支持度和置信度了，其中支持度  $support = \frac{rulesupCount}{|D|} \times 100\%$ ，置信度  $confidence = \frac{rulesupCount}{condsupCount} \times 100\%$ 。那么对于那个例子，我们不难计算出  $support = \frac{2}{10} \times 100\% = 20\%$ ， $confidence = \frac{2}{3} \times 100\% = 66.7\%$ 。

如果一个规则项集满足最小支持度约束  $minsup$ ，则其是一个频繁规则项集，与之相对应的则是非频繁的规则项集。依然看前面的例子，假如我们规定  $minsup = 10\%$ ，那么  $support = 20\% > minsup = 10\%$ ，因而这条规则是频繁的。如果一个规则的置信度满足最小置信度约束  $minconf$ ，那么它就是一个准确的规则。类关联规则集合中的规则应当都是频繁且准确的规则。

### 2.2.2 规则生成的 CBA-RG 算法

文献 [1] 所提出的 CBA-RG 算法事实上是在 Apriori 算法的基础进行修改得到的，因而其基本思想与 Apriori 算法十分接近。为了便于我们后面的讨论，这里将 [1] 中的算法抄录于此，见算法 2。在这个算法中，我们将一个规则项集表示为形如

$$< (condset, condsupCount), (y, rulesupCount) >$$

且称  $condset$  内拥有  $k$  项元素的规则项集为  $k$ -ruleitems。

---

#### 算法 2 CBA-RG 算法

---

```

1:  $F_1 = \{\text{large 1-ruleitems}\}$ 
2:  $CAR_1 = \text{genRules}(F_1)$ 
3:  $prCAR_1 = \text{pruneRules}(CAR_1)$ 
4:  $k = 2$ 
5: while  $F_{k-1} \neq \emptyset$  do
6:    $C_k = \text{candidateGen}(F_{k-1})$ 
7:   for each data case  $d \in D$  do
8:      $C_d = \text{ruleSubset}(C_k, d)$ 
9:     for each candidate  $c \in C_d$  do
10:       $c.condsupCount++$ 
11:      if  $d.class = c.class$  then
12:         $c.rulesupCount++$ 
13:      end if
14:    end for
15:  end for
16:   $F_k = \{c \in C_k | c.rulesupCount \leq minsup\}$ 
17:   $CAR_k = \text{genRules}(F_k)$ 
18:   $prCAR_k = \text{pruneRules}(CAR_k)$ 
19:   $k++$ 
20: end while
21:  $CARs = \bigcup_k CAR_k$ 
22:  $prCARs = \bigcup_k prCAR_k$ 

```

---

刘兵等人已在 [1] 中对上述算法进行了详细的讨论，我们不再这里做重复的分析，我们下面只分析作者没有在论文中指明的步骤。

一个问题就是算法 2 的第 2 行 **genRules** 的实现。我们在 2.2.1 节中给出了如何从规则项集生成类关联规则的过程。但是这个当面对一个规则项集时,简单的采用这个方法会存在一些问题。假如存在两个规则项集,它们拥有相同的 *condset*, 那么我们应当选取那个具有最高置信度的规则项集。如果这两个规则项集的置信度也一致,那么就随机取一条作为该 *condset* 的代表规则项集。依然拿前面表 4 为例,并假设  $minconf = 60\%$ , 我们可以找到 2 个 *condset* 一致的 2-ruleitems:

$$< (\{(A, g), (B, q)\}, 3), ((C, y), 2) >$$

$$< (\{(A, g), (B, q)\}, 3), ((C, n), 1) >$$

这两个规则项集的 *condset* 一致,只是最后的所属类不同。对于第一个规则项集,其置信度为  $\frac{2}{3}$ , 第二个规则项集的置信度为  $\frac{1}{3}$ 。显然前者的置信度更高,所以我们选取第一个规则项集以生成类关联规则

$$\{(A, g), (B, q)\} \rightarrow (C, y)$$

因此,在实现 **genRules** 时,需要按照上述策略处理 *condset* 一致的情形。详细的算法描述如算法 3 所示。

---

### 算法 3 **genRules** 的实现

---

```

1: function genRules(F)
2:   CAR =  $\emptyset$ 
3:   for each ruleitem r  $\in$  F do
4:     if r.condset in CAR then
5:       if r.confidence >  $\max(\{c.confidence : c \in CAR \wedge c.condset = r.condset\})$  then
6:         replace c with r
7:       end if
8:     else
9:       CAR = CAR  $\cup$  r
10:    end if
11:  end for
12:  return CAR
13: end function

```

---

第二个需要解决的问题就是规则剪枝的问题,即 **pruneRules** 的实现。文献 [1] 中作者介绍到其采用的是 C4.5 算法中的利用悲观错误率进行剪枝的方法。给定一个类关联规则  $r : A \rightarrow y$ , 考虑剪枝后的规则  $r' : A' \rightarrow y$ , 其中  $A'$  是  $A$  中去掉一个合取项后得到的。只要剪枝后的规则的悲观误差率不高于原规则的误差率,就保留其中悲观误差率最低的规则。重复规则剪枝步骤,直到规则的悲观误差率不能再进行改进为止。由于某些规则在剪枝后会变得相同,因而需要丢弃重复的规则。在上面的叙述中,我们需要强调剪枝的条件是剪枝后的规则的悲观误差率不高于原规则的误差率,这里的“不高于”而非严格的“低于”是基于众所周知的奥卡姆剃刀,即

**奥卡姆剃刀:** 给定两个具有相同泛化误差的模型,较简单的模型比较复杂的模型更可取。

这个算法的形式化的描述见算法 4。

第三个需要解决的问题就是如何实现 **candidateGen**。**candidateGen** 函数和 Apriori 算法中的 **Apriori-gen** 函数十分类似。该操作由前一次迭代发现的  $F_{k-1}$  产生新的候选  $k$ -ruleitems。我们这里采用  $F_{k-1} \times F_{k-1}$  方法。由于 **candidateGen** 函数合并一对频繁  $(k-1)$ -ruleitems,当且仅当它们的前  $k-2$  个项都相同。但是这么做还不能满足先验原理,因为合并后的项的子集不一定是频繁的,因而还需要进行剪枝,即考虑候选



---

**算法 4** pruneRules 的实现

---

```
1: function pruneRules( $r$ )
2:   if  $|r.condset| = 0$  then return
3:   else
4:      $prs = \emptyset$ 
5:      $r' = r$ 
6:     for each item  $t$  in  $r.condset$  do
7:        $r'.condset = r.condset - t$ 
8:        $prs = prs \cup r'$ 
9:     end for
10:     $pr = \arg \min \{r'.errorRate : r' \in prs\}$ 
11:    if  $pr.errorRate \leq r.errorRate$  then
12:      return pruneRules( $pr$ )
13:    end if
14:  end if
15: end function
```

---

$k$ -ruleitems  $X = \{i_1, i_2, \dots, i_k\}$ , 如果存在它的一个真子集  $X - \{i_j\} (\forall j = 1, 2, \dots, k)$  是非频繁的, 那么  $X$  将被剪枝。用形式化的语言描述 [11], 有

$$\text{candidateGen}(F_{k-1}) = \{a \cup \{b\} | a \in F_{k-1} \wedge b \notin a\} - \{c | \{s | s \subseteq c \wedge |s| = k-1\} \not\subseteq F_{k-1}\}$$

至此, 我们就不难实现 CBA-RG 算法了。

## 2.3 构建分类器的 CBA-CB M1 算法实现

### 2.3.1 相关概念

我们称规则  $r : condset \rightarrow y$  覆盖样本  $d$ , 当且仅当  $r$  的前件  $condset$  满足样本  $d$ , 即对应的属性值完全相同。如果这条规则的后件  $y$  与样本  $d$  的分类一致, 则称这条规则正确分类了样本  $d$ 。

下面定义类关联规则的偏序。给定两个规则  $r_i, r_j$ , 定义  $r_i > r_j$  当

- a)  $r_i$  的置信度高于  $r_j$  的置信度;
- b)  $r_i$  和  $r_j$  的置信度一致, 但是  $r_i$  的支持度高于  $r_j$  的支持度;
- c)  $r_i$  和  $r_j$  的置信度和支持度都一致, 但是  $r_i$  比  $r_j$  的生成时间要早。

CBA-CB 算法的基本思想就是构建一个有序规则序列

$$\langle r_1, r_2, \dots, r_n, default\_class \rangle$$

作为最后的分类器。其中  $r_i \in R$ ,  $R$  为生成的类关联规则集合。序列中的两条规则  $r_a, r_b$ , 当  $b > a$  时有  $r_a > r_b$ 。  $default\_class$  是默认类, 当前面的那些规则都不能覆盖这条记录时, 则将其划分入默认类中。

### 2.3.2 CBA-CB M1 算法的实现

这个略显幼稚的算法的主要思想是基于贪心策略, 尽可能早的选出高秩的那些规则。为了更好的说明这个算法的思想, 我们用图示进行解释。图 1(a)展示了初始时数据集的情况, 其中包括若干正例和反例, 虚

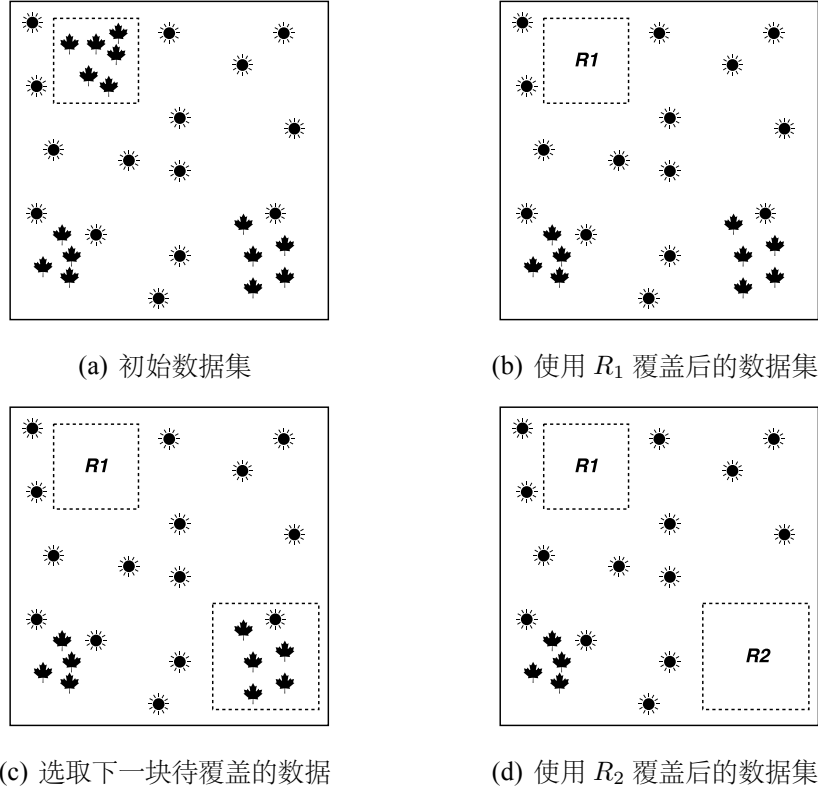


图 1 CBA-CB M1 算法示例

线框出区域内是秩最高的规则  $R_1$  所能覆盖的样本。去除了  $R_1$  所能覆盖的所有训练记录后，剩下的数据样本如图 1(b)所示。接着在剩余的样本中寻找秩次高的规则  $R_2$  所能覆盖的样本，如图 1(c)中的虚线框所示。然后删去  $R_2$  所能覆盖的这些样本，留下的记录如图 1(d)所示。算法描述见算法 5所示。

## 2.4 构建分类器的 CBA-CB M2 算法实现

尽管 CBA-CB M1 算法形式相对简单，易于编程实现，但是当面对大训练样本的情形时，其效率不高。所以文献 [1] 中又紧接着提出了 CBA-CB M2 算法，其中心思想同 M1 基本一致。该算法包括 3 个阶段，通过首先缩小候选规则的方法，减少不必要的规则判定，提高训练效率。算法的详细解释可见文献 [1]，这里只提供伪代码描述，见算法 6、7和 8。

但是需要注意的是，我们深切怀疑算法 8第 12 行与第 13 行之间漏写了一步操作，这里应当将规则  $r$  所覆盖的那些训练记录从样本集中删除，否则后续第 14 至 15 行更新当前数据集类分布的工作就无需进行了，因为样本集根本没有发生变化。不论是 M1 还是 M2 算法，其大体的框架应该是类似的，否则两个版本的算法得到的分类器完全不同，这显然不符合常理。

在实际编写程序时，由于某些数据集所包含的规则太多，内存无法存放下，而且十分耗费时间。因而我们考虑当得到的类关联规则数达到 2000，或者收敛速度已呈现亚线性（即每一轮迭代所挖掘出的规则数少于 10 个）时，则不再继续进行挖掘，提前终止迭代。实验结果表明，采取这样的策略对于最后分类器的性能并没有太大影响。我们认为，这是由于高秩规则已经可以很好的描述数据集中存在的联系，过多的规则对于提高分类器性能的影响有限。

---

**算法 5 CBA-CB M1 算法**

---

```
1:  $R = \text{sort}(R)$ 
2: for each rule  $r \in R$  in sequence do
3:    $temp = \emptyset$ 
4:   for each case  $d \in D$  do
5:     if  $d$  satisfies the conditions of  $r$  then
6:       store  $d.id$  in  $temp$  and mark  $r$  if it correctly classifies  $d$ 
7:     end if
8:   end for
9:   if  $r$  is marked then
10:    insert  $r$  at the end of  $C$ 
11:    delete all the cases with the ids in  $temp$  from  $D$ 
12:    selecting a default class for the current  $C$ 
13:    compute the total number of errors of  $C$ 
14:   end if
15: end for
16: Find the first rule  $p$  in  $C$  with the lowest total number of errors and drop all the rules after  $p$  in  $C$ 
17: Add the default class associated with  $p$  to end of  $C$ , and return  $C$  (our classifier)
```

---

---

**算法 6 CBA-CB M2 算法: Stage 1**

---

```
1:  $Q = \emptyset; U = \emptyset; A = \emptyset$ 
2: for each case  $d \in D$  do
3:    $cRule = \text{maxCoverRule}(C_c, d)$ 
4:    $wRule = \text{maxCoverRule}(C_w, d)$ 
5:    $U = U \cup \{cRule\}$ 
6:    $cRule.classCasesCovered[d.class]++$ 
7:   if  $cRule > wRule$  then
8:      $Q = Q \cup \{cRule\}$ 
9:     mark  $cRule$ 
10:  else  $A = A \cup \langle d.id, d.class, cRule, wRule \rangle$ 
11:  end if
12: end for
```

---

---

**算法 7 CBA-CB M2 算法: Stage 2**

---

```
1: for each entry  $\langle dID, y, cRule, wRule \rangle \in A$  do
2:   if  $wRule$  is marked then
3:      $cRule.classCasesCovered[y] --$ 
4:      $wRule.classCasesCovered[y] ++$ 
5:   else
6:      $wSet = allCoverRules(U, dID.case, cRule)$ 
7:     for each rule  $w \in wSet$  do
8:        $w.replace = w.replace \cup \{\langle cRule, dID, y \rangle\}$ 
9:        $w.classCasesCovered[y] ++$ 
10:    end for
11:     $Q = Q \cup wSet$ 
12:   end if
13: end for
```

---

---

**算法 8 CBA-CB M2 算法: Stage 3**

---

```
1:  $classDistr = compClassDistr(D)$ 
2:  $ruleErrors = 0$ 
3:  $Q = sort(Q)$ 
4: for each rule  $r$  in  $Q$  in sequence do
5:   if  $r.classCasesCovered[r.class] \neq 0$  then
6:     for each entry  $\langle rul, dID, y \rangle$  in  $r.replace$  do
7:       if the  $dID$  case has been covered by a previous  $r$  then
8:          $r.classCasesCovered[y] --$ 
9:       else
10:         $rul.classCasesCovered[y] --$ 
11:      end if
12:    end for
13:     $ruleErrors = ruleErrors + errorsOfRule(r)$ 
14:     $classDistr = update(r, classDistr)$ 
15:     $defaultClass = selectDefault(classDistr)$ 
16:     $totalErrors = ruleErrors + defaultErrors$ 
17:    Insert  $\langle r, default-class, totalErrors \rangle$  at the end of  $C$ 
18:   end if
19: end for
20: Find the first rule  $p$  in  $C$  with the lowest  $totalErrors$ , and then discard all the rules after  $p$  from  $C$ 
21: Add the default class associated with  $p$  to end of  $C$ 
22: Return  $C$  without  $totalErrors$  and  $default-class$ 
```

---

### 3 实验结果与分析

编写程序重现文献 [1] 的工作是本文的中心工作，因而如何检验我们的实现的程序至关重要。我们仿照文章中的检验手段，采用了 10 折交叉检验去计算分类器的错误率。因而在正式给出我们的测试结果之前，首先介绍一下 10 折交叉检验。

#### 3.1 10 折交叉检验

所谓 10 折交叉检验 [8]，就是将数据随机分成大小相同的 10 份，在每次运行时，选择其中一份作为检验集，而其余的全部作为训练集，该过程重复进行 10 次，使得每份数据都用于检验恰好一次。最后取 10 次运行的误差的平均值（一说是总和 [12]）作为最后的误差估计。限于篇幅，图 2 给出了 3 折交叉检验的简单示意，10 折交叉检验的思想和其基本一致。

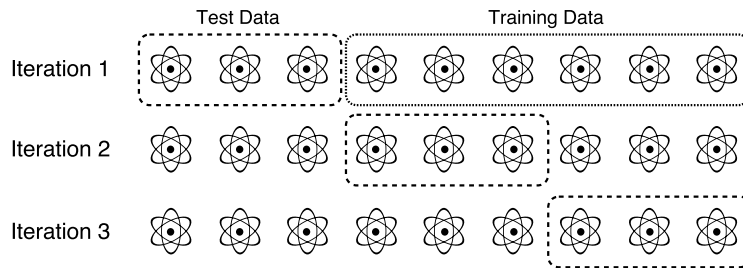


图 2 3 折交叉检验的简单示意图

#### 3.2 样例测试

本文的实现程序全部使用 Python 3.6.0 编写，运行在 2.9 GHz Intel Core i5 处理器上。我们从 UCI 机器学习资源库 [2] 中选取了 30 个数据集，对我们的程序进行测试。我们的运行时间全部基于驻留在内存的数据集。所有指标都采用 10 折交叉检验，取 10 次检验的平均性能以评估其总体性能指标。具体的测试结果如表 5 和表 6 所示。表中每一行对应一个数据集，限于篇幅，数据集名称可能存在简写。表格分成两大大部分，左侧 3 至 8 列对应于不进行剪枝时的各种测试结果，右侧 9 至 13 列对应剪枝下的测试结果。以左侧为例，第 3 列表示 10 折交叉检验下的平均错误率；第 4 列为 CBA-RG 算法生成的类关联规则的平均数量，若考虑剪枝，则这一列表示的是剪枝后的规则数；第 5 列记录 CBA-RG 算法的平均执行时间，若考虑剪枝，则这一部分时间还涵盖了剪枝所花费的时间；第 6 列表示 CBA-CB 算法的平均执行时间；第 7 列表示最后训练得到的分类器中所包含的类关联规则的平均数量。

从测试结果不难看出，我们训练出的分类器，其判定错误率较低，基本可以满足分类任务的要求。但是必须看到，我们的程序，尤其是 CBA-RG 算法的运行时间过长。我们分析，认为有以下几点原因：a) 我们所采用的 Python 语言是一种动态语言，其执行效率相比于论文中采用的 C++ 语言肯定要低许多；b) 为了调试程序的方便，我们在程序中额外增加了许多事实上与类关联规则和分类器无关的其他信息，这势必会增加程序运行的开销，降低程序的运行效率；c) 我们所采用的数据结构较为简陋，理论上应当采用树型结构实现规则和分类器。数据结构对算法执行效率的影响是主要的。

综上，我们可以总结，我们实现的基于类关联规则的分类器，其判定准确率较高，但是执行时间较慢，有待以后的优化提高。

表 5 使用 CBA-CB M1 方法下的测试结果

序号	数据集	不剪枝				剪枝						
		错误率/%	#CARS	RG 时间/s	CB 时间/s	# 分类器	CARS	错误率/%	#CARS	RG 时间/s	CB 时间/s	# 分类器
1	australian	1.9	1140	15.48	0.19	58	752	2.9	752	75.46	0.22	70
2	banknote	0.1	194	0.90	0.09	28	68	0.0	68	2.17	0.06	28
3	car	1.3	254	2.30	0.08	24	199	1.1	199	3.45	0.08	26
4	diagnosis	0.0	692	2.97	0.00	4	126	0.0	126	5.48	0.00	4
5	dresses	3.0	511	2.64	0.28	89	414	2.0	414	2.45	0.24	89
6	facebook	0.8	1995	8.19	0.15	59	1254	1.2	1254	70.19	0.17	61
7	flare	1.0	567	6.17	0.04	44	423	1.1	423	20.05	0.04	46
8	forest	3.4	2116	2.88	0.18	40	1658	3.5	1658	26.34	0.16	40
9	german	0.0	1294	10.43	0.48	85	906	0.1	906	15.74	0.30	80
10	heart-disease	1.6	659	16.75	0.06	35	475	2.4	475	39.65	0.05	34
11	horse-colic	0.0	2068	13.65	0.14	57	1594	0.7	1594	44.12	0.17	38
12	ilpd	0.2	1002	18.48	0.22	64	666	0.4	666	21.42	0.21	67
13	ionosphere	1.1	2008	12.56	0.13	41	1378	3.5	1378	34.76	0.11	46
14	iris	0.7	120	0.09	0.00	9	45	0.7	45	0.20	0.00	9
15	kr-vs-kp	0.0	2241	11.32	0.95	39	1149	0.1	1149	69.80	0.78	29
16	led	5.0	2	0.01	0.00	2	1	17.5	1	0.01	0.00	1
17	mammo	0.0	306	3.04	0.24	64	220	0.5	220	4.02	0.19	65
18	messidor	9.6	2359	9.04	1.27	86	1879	9.9	1879	30.59	1.12	73
19	monks	0.0	656	2.27	0.04	40	455	0.0	455	4.53	0.01	22
20	pima	2.3	751	14.58	0.24	68	630	2.4	630	33.56	0.19	52
21	quali-bank	0.0	731	2.89	0.01	9	141	0.0	141	10.27	0.00	10
22	seeds	1.9	749	3.87	0.01	23	168	1.9	168	17.79	0.01	23
23	tic-tac-toe	0.1	600	2.15	0.21	61	540	0.1	540	3.33	0.13	55
24	transfusion	4.3	9	0.00	0.00	4	9	2.1	9	0.01	0.00	4
25	user-know	4.3	7	0.00	0.00	5	7	4.3	7	0.01	0.00	5
26	vertebra	0.0	521	8.67	0.04	15	133	0.0	133	12.81	0.01	15
27	wiki4HE	0.1	1150	5.36	0.42	86	720	0.3	720	45.12	0.33	79
28	wine	0.6	1702	26.28	0.05	28	1004	1.1	1004	83.82	0.04	27
29	yeast	3.1	3	0.11	0.01	3	3	3.1	3	0.12	0.02	3
30	zoo	1.0	2397	5.95	0.03	10	1780	1.0	1780	38.76	0.02	11

表 6 使用 CBA-CB M2 方法下的测试结果

序号	数据集	不剪枝				剪枝							
		错误率/%	#CARS	RG 时间/s	CB 时间/s	# 分类器	CARS	错误率/%	#CARS	RG 时间/s	CB 时间/s	# 分类器	CARS
1	australian	0.0	1225	11.69	0.63	15	905	0.0	905	99.24	0.21	76	76
2	banknote	0.0	194	0.86	0.24	13	67	0.1	67	2.18	0.06	27	27
3	car	0.6	276	2.64	0.30	8	192	1.6	192	4.14	0.07	23	23
4	diagnosis	0.0	695	2.89	0.08	4	125	0.0	125	5.73	0.00	4	4
5	dresses	0.0	536	2.29	0.20	29	410	2.2	410	2.31	0.20	85	85
6	facebook	1.0	1825	26.40	0.63	10	959	1.8	959	90.78	0.17	58	58
7	flare	0.0	475	57.25	0.17	19	379	0.3	379	99.89	0.17	20	20
8	forest	1.5	2285	5.22	0.82	12	1645	2.1	1645	30.43	0.76	13	13
9	german	0.2	1281	12.44	0.86	13	884	0.2	884	31.78	0.73	9	9
10	heart-disease	1.0	492	9.52	0.14	6	374	1.9	374	37.69	0.12	10	10
11	horse-colic	0.4	1978	31.44	0.61	32	1238	0.9	1238	83.42	0.52	29	29
12	ilpd	0.3	1015	20.00	0.43	1	780	0.4	780	53.17	0.37	4	4
13	ionosphere	0.0	1975	6.86	0.68	18	1499	0.3	1499	30.12	0.56	20	20
14	iris	1.3	120	0.07	0.02	7	44	0.0	44	0.20	0.00	8	8
15	kr-vs-kp	0.0	2202	9.60	5.53	12	1308	0.1	1308	70.34	4.67	21	21
16	led	10.0	2	0.01	0.00	1	1	15.0	1	0.01	0.00	1	1
17	mammo	0.0	299	2.71	0.25	8	220	0.3	220	3.93	0.20	64	64
18	messidor	0.1	1389	13.29	1.52	3	1123	0.3	1123	40.12	1.31	5	5
19	monks	0.0	683	3.05	0.20	34	415	0.0	415	4.30	0.01	21	21
20	pima	0.0	740	11.98	0.43	6	681	0.3	681	42.70	0.39	8	8
21	quali-bank	0.0	748	3.24	0.16	7	147	0.0	147	10.44	0.00	10	10
22	seeds	1.0	790	3.93	0.15	20	170	2.4	170	18.99	0.00	23	23
23	tic-tac-toe	0.0	604	2.19	0.38	19	542	0.0	542	3.52	0.14	56	56
24	transfusion	0.0	9	0.00	0.01	2	9	0.0	9	0.01	0.00	4	4
25	user-know	5.2	7	0.00	0.00	3	7	4.4	7	0.01	0.00	5	5
26	vertebra	0.0	523	8.44	0.16	3	132	0.0	132	13.24	0.01	15	15
27	wiki4HE	0.9	1064	4.51	0.64	1	789	1.1	789	30.65	0.49	5	5
28	wine	0.1	1670	52.63	0.30	29	1232	0.2	1232	97.79	0.25	27	27
29	yeast	3.1	3	0.13	0.02	2	3	3.1	3	0.12	0.01	3	3
30	zoo	1.0	2325	8.78	0.28	8	1892	1.1	1892	53.20	0.24	10	10

## 4 致谢与后记

数据挖掘与知识发现的课程项目至此快要进入尾声了。从三月份完成选题、编制项目计划书，到现在完成了这份报告的编写，一路的酸甜苦辣已回荡在脑海之中。

陈大松同学完成了数据预处理的工作，编写了缺失值处理和连续型属性值离散化的相关程序，为我们的工作开了个好头。刘砺志同学主要负责 CBA 算法的实现工作，并且完成了这份报告的编写任务，并将其交由另外两名同学审阅并进行修改。肖璐菁同学负责数据集的格式归一化处理，并编写了测试程序，同时记录了各项数据指标。这项工作十分枯燥无味，但是肖同学还是坚持并优秀的完成了它们。感谢大家团结一致、戮力同心，坚持不懈的完成了论文的重现工作。

最后还要感谢林琛老师的博学多才，为我们开启了数据挖掘这个新世界的大门，带领我们领略数据科学的独特魅力。通过这项工作，让我们体会到数据科学的乐趣，为今后的从事相关工作奠定了坚实的基础。

## 参考文献

- [1] Liu B, Hsu W, Ma Y. Integrating Classification and Association Rule Mining. Proc of Kdd, 1998:80–86.
- [2] UC Irvine Machine Learning Repository, <http://archive.ics.uci.edu/ml/>
- [3] Rakesh Agrawal, Ramakrishnan Srikant. Fast algorithms for mining association rules in large databases. Proceedings of the 20th International Conference on Very Large Data Bases, VLDB, pages 487–499, Santiago, Chile, September 1994.
- [4] Ho, Tin Kam. Random Decision Forests. Proceedings of the 3rd International Conference on Document Analysis and Recognition, Montreal, QC, 14–16 August 1995. pp. 278–282.
- [5] Catlett J. On changing continuous attributes into ordered discrete attributes. European Working Session on Learning on Machine Learning. Springer-Verlag New York, Inc. 1991:164–178.
- [6] Fayyad U M. Multi-Interval Discretization of Continuous-Valued Attributes for Classification Learning. International Joint Conference on Artificial Intelligence. 1993:1022–1027.
- [7] Dougherty J, Kohavi R, Sahami M. Supervised and Unsupervised Discretization of Continuous Features. Machine Learning Proceedings, 1995(2):194–202.
- [8] Cross-validation (statistics), [https://en.wikipedia.org/wiki/Cross-validation\\_\(statistics\)](https://en.wikipedia.org/wiki/Cross-validation_(statistics))
- [9] Discretizing a continuous variable using Entropy, <http://www.clear-lines.com/blog/post/Discretizing-a-continuous-variable-using-Entropy.aspx>
- [10] Integrating Classification and Association Rule Mining - the Secret Behind CBA, <https://www.cs.uic.edu/~hxiao/courses/cs594-slides.pdf>
- [11] Apriori algorithm, [https://en.wikipedia.org/wiki/Apriori\\_algorithm](https://en.wikipedia.org/wiki/Apriori_algorithm)
- [12] Tan P N, Steinbach M, Kumar V. Introduction to Data Mining, (First Edition). Addison-Wesley Longman Publishing Co. Inc. 2005.



## 附录 A 相关源程序

本文使用的程序均采用 Python 3.6.0 编写，运行在安装 Mac OS 10.12.2 操作系统、搭载有 2.9 GHz Intel Core i5 处理器、使用 8 GB 2133 MHz LPDDR3 内存的计算机上。

### 一、程序清单rmep.py

```
1  """
2  Description: Recursive minimal entropy partitioning, to discretize
              continuous-valued attributes. We use the supervised
3      algorithm presented in Fayyad & Irani (1993) and introduced in
              Dougherty, Kohavi & Sahami (1995) section 3.3.
4      We also refer to a F# code on GitHub (https://gist.github.com/mathias-brandewinder/5650553).
5  Input: a data table with several rows but only two column, the first
              column is continuous-valued (numerical) attributes,
6      and the second column is the class label of each data case (
              categorical).
7      e.g. data = [[1.0, 'Yes'], [0.5, 'No'], [2.0, 'Yes']]
8  Output: a list of partition boundaries of the range of continuous-valued
              attribute in ascending sort order.
9      e.g. walls = [0.5, 0.8, 1.0], thus we can separate the range into 4
              intervals: <=0.5, 0.5<*<=0.8, 0.8<*<=1.0 & >=1.0
10 Author: CBA Studio
11 Reference:
12     1. Multi-Interval Discretization of Continuous-Valued Attributes for
              Classification Learning, Fayyad & Irani, 1993
13     2. Supervised and Unsupervised Discretization of Continuous Features,
              Dougherty, Kohavi & Sahami, 1995
14     3. http://www.clear-lines.com/blog/post/Discretizing-a-continuous-
              variable-using-Entropy.aspx
15 """
16 import math
17
18
19 # A block to be split
20 # It has 4 member:
21 #   data: the data table with a column of continuous-valued attribute and
              a column of class label
22 #   size: number of data case in this table
23 #   number_of_classes: obviously, the number of class in this table
24 #   entropy: entropy of dataset
25 class Block:
26     def __init__(self, data):
27         self.data = data
28         self.size = len(data)
29         classes = set([x[1] for x in data])      # get distinct class
              labels in this table
30         self.number_of_classes = len(set(classes))
```

```

31         self.entropy = calculate_entropy(data)
32
33
34 # Calculate the entropy of dataset
35 # parameter data: the data table to be used
36 def calculate_entropy(data):
37     number_of_data = len(data)
38     classes = set([x[1] for x in data])
39     class_count = dict([(label, 0) for label in classes])
40     for data_case in data:
41         class_count[data_case[1]] += 1      # count the number of data
42                                             # case of each class
43     entropy = 0
44     for c in classes:
45         p = class_count[c] / number_of_data
46         entropy -= p * math.log2(p)          # calculate information
47                                             # entropy by its formula, where the base is 2
48
49     return entropy
50
51 # Compute Gain(A, T: S) mentioned in Dougherty, Kohavi & Sahami (1995), i
52 # .e. entropy gained by splitting original_block
53 # into left_block and right_block
54 # original_block: the block before partition
55 # left_block: the block split which its value below boundary
56 # right_block: the block above boundary
57 def entropy_gain(original_block, left_block, right_block):
58     gain = original_block.entropy - \
59         ((left_block.size / original_block.size) * left_block.entropy
60          +
61          (right_block.size / original_block.size) * right_block.
62          entropy)
63     return gain
64
65 # Get minimum entropy gain required for a split of original_block into 2
66 # blocks "left" and "right", see Dougherty,
67 # Kohavi & Sahami (1995)
68 # original_block: the block before partition
69 # left_block: the block split which its value below boundary
70 # right_block: the block above boundary
71 def min_gain(original_block, left_block, right_block):
72     delta = math.log2(math.pow(3, original_block.number_of_classes) - 2)
73     - \
74         (original_block.number_of_classes * original_block.entropy -
75          left_block.number_of_classes * left_block.entropy -
76          right_block.number_of_classes * right_block.entropy)
77     gain_sup = math.log2(original_block.size - 1) / original_block.size +

```

```

        delta / original_block.size
72     return gain_sup
73
74
75 # Identify the best acceptable value to split block
76 # block: a block of dataset
77 # Return value: a list of (boundary, entropy gain, left block, right
    block) or
78 #     None when it's unnecessary to split
79 def split(block):
80     candidates = [x[0] for x in block.data]      # candidates is a list of
        values can be picked up as boundary
81     candidates = list(set(candidates))          # get different values in
        table
82     candidates.sort()                          # sort ascending
83     candidates = candidates[1:]                # discard smallest,
        because by definition no value is smaller
84
85     wall = []      # wall is a list storing final boundary
86     for value in candidates:
87         # split by value into 2 groups, below & above
88         left_data = []
89         right_data = []
90         for data_case in block.data:
91             if data_case[0] < value:
92                 left_data.append(data_case)
93             else:
94                 right_data.append(data_case)
95
96         left_block = Block(left_data)
97         right_block = Block(right_data)
98
99         gain = entropy_gain(block, left_block, right_block)
100        threshold = min_gain(block, left_block, right_block)
101
102        # minimum threshold is met, the value is an acceptable candidate
103        if gain >= threshold:
104            wall.append([value, gain, left_block, right_block])
105
106    if wall:      # has candidate
107        wall.sort(key=lambda wall: wall[1], reverse=True)  # sort
            descending by "gain"
108        return wall[0]      # return best candidate with max entropy gain
109    else:
110        return None      # no need to split
111
112
113 # Top-down recursive partition of a data block, append boundary into "

```

```

        walls"
114 # block: a data block
115 def partition(block):
116     walls = []
117
118     # inner recursive function, accumulate the partitioning values
119     # sub_block: just a data block
120     def recursive_split(sub_block):
121         wall_returned = split(sub_block)           # binary partition, get
                                                    # bin boundary
122         if wall_returned:                           # still can be spilt
123             walls.append(wall_returned[0])          # record this
                                                    # partitioning value
124             recursive_split(wall_returned[2])       # recursively process
                                                    # left block
125             recursive_split(wall_returned[3])       # recursively split right
                                                    # block
126         else:
127             return                                  # end of recursion
128
129     recursive_split(block)                          # call inner function
130     walls.sort()                                    # sort boundaries descending
131     return walls
132
133
134 # just for test
135 if __name__ == '__main__':
136     import random
137
138     test_data = []
139     for i in range(100):
140         test_data.append([random.random(), random.choice(range(0, 2))])
141         test_data.append([random.random() + 1, random.choice(range(2, 4))
142                             ])
143         test_data.append([random.random() + 2, random.choice(range(4, 6))
144                             ])
145         test_data.append([random.random() + 3, random.choice(range(6, 8))
146                             ])
147
148     test_block = Block(test_data)
149     test_walls = partition(test_block)
150     print(test_walls)          # should be [1+e, 2+e, 3+e], where e is a
                                # number very close to 0

```

## 二、程序清单pre\_processing.py

```

1 """
2 Description: Pre-process original data. Firstly, we process the missing
  values (donated as '?'), discarding this column

```

```

3     when missing ratio above 50%, or filling blanks when below. We "guess
      " missing values by simply filling the mode of
4     existing values in the same column. And then, for the numerical
      attribute, we discretize it by recursive minimal
5     entropy partitioning (see rmep.py). For the categorical attribute, we
      just replace the label with a
6     positive integer. For more information, see [1].
7 Input: a data table with several data case, many attributes and class
      label in the last column, a list of the name of
8     each attribute, and a list of the type of each column.
9 Output: a data list without numerical values and "str" categorical values
      .
10 Author: CBA Studio
11 Reference:
12     1. http://cgi.csc.liv.ac.uk/~frans/KDD/Software/LUCS-KDD-DN/lucs-
      kdd\_DN.html
13 """
14 import rmep
15
16
17 # Identify the mode of a list, both effective for numerical and
      categorical list. When there exists too many modes
18 #   having the same frequency, return the first one.
19 # arr: a list need to find mode
20 def get_mode(arr):
21     mode = []
22     arr_appear = dict((a, arr.count(a)) for a in arr)    # count
      appearance times of each key
23     if max(arr_appear.values()) == 1:                    # if max time is 1
24         return                                           # no mode here
25     else:
26         for k, v in arr_appear.items():                  # else, mode is the number
      which has max time
27             if v == max(arr_appear.values()):
28                 mode.append(k)
29     return mode[0]   # return first number if has many modes
30
31
32 # Fill missing values in column column_no, when missing values ration
      below 50%.
33 # data: original data list
34 # column_no: identify the column No. of that to be filled
35 def fill_missing_values(data, column_no):
36     size = len(data)
37     column_data = [x[column_no] for x in data]           # get that column
38     while '?' in column_data:
39         column_data.remove('?')
40     mode = get_mode(column_data)

```

```

41     for i in range(size):
42         if data[i][column_no] == '?':
43             data[i][column_no] = mode # fill in mode
44     return data
45
46
47 # Get the list needed by rmep.py, just glue the data column with class
   column.
48 # data_column: the data column
49 # class_column: the class label column
50 def get_discretization_data(data_column, class_column):
51     size = len(data_column)
52     result_list = []
53     for i in range(size):
54         result_list.append([data_column[i], class_column[i]])
55     return result_list
56
57
58 # Replace numerical data with the No. of interval, i.e. consecutive
   positive integers.
59 # data: original data table
60 # column_no: the column No. of that column
61 # walls: the split point of the whole range
62 def replace_numerical(data, column_no, walls):
63     size = len(data)
64     num_spilt_point = len(walls)
65     for i in range(size):
66         if data[i][column_no] > walls[num_spilt_point - 1]:
67             data[i][column_no] = num_spilt_point + 1
68             continue
69         for j in range(0, num_spilt_point):
70             if data[i][column_no] <= walls[j]:
71                 data[i][column_no] = j + 1
72                 break
73     return data
74
75
76 # Replace categorical values with a positive integer.
77 # data: original data table
78 # column_no: identify which column to be processed
79 def replace_categorical(data, column_no):
80     size = len(data)
81     classes = set([x[column_no] for x in data])
82     classes_no = dict([(label, 0) for label in classes])
83     j = 1
84     for i in classes:
85         classes_no[i] = j
86         j += 1

```

```

87     for i in range(size):
88         data[i][column_no] = classes_no[data[i][column_no]]
89     return data, classes_no
90
91
92 # Discard all the column with its column_no in discard_list
93 # data: original data set
94 # discard_list: a list of column No. of the columns to be discarded
95 def discard(data, discard_list):
96     size = len(data)
97     length = len(data[0])
98     data_result = []
99     for i in range(size):
100         data_result.append([])
101         for j in range(length):
102             if j not in discard_list:
103                 data_result[i].append(data[i][j])
104     return data_result
105
106
107 # Main method here, see Description in detail
108 # data: original data table
109 # attribute: a list of the name of attribute
110 # value_type: a list identifying the type of each column
111 # Returned value: a data table after process
112 def pre_process(data, attribute, value_type):
113     column_num = len(data[0])
114     size = len(data)
115     class_column = [x[-1] for x in data]
116     discard_list = []
117     for i in range(0, column_num - 1):
118         data_column = [x[i] for x in data]
119
120         # process missing values
121         missing_values_ratio = data_column.count('?') / size
122         if missing_values_ratio > 0.5:
123             discard_list.append(i)
124             continue
125         elif missing_values_ratio > 0:
126             data = fill_missing_values(data, i)
127             data_column = [x[i] for x in data]
128
129         # discretization
130         if value_type[i] == 'numerical':
131             discretization_data = get_discretization_data(data_column,
132                                                             class_column)
133             block = rmep.Block(discretization_data)
134             walls = rmep.partition(block)

```

```

134         if len(walls) == 0:
135             max_value = max(data_column)
136             min_value = min(data_column)
137             step = (max_value - min_value) / 3
138             walls.append(min_value + step)
139             walls.append(min_value + 2 * step)
140             print(attribute[i] + ":", walls)          # print out split
                                                    points
141             data = replace_numerical(data, i, walls)
142         elif value_type[i] == 'categorical':
143             data, classes_no = replace_categorical(data, i)
144             print(attribute[i] + ":", classes_no)      # print out
                                                    replacement list
145
146         # discard
147         if len(discard_list) > 0:
148             data = discard(data, discard_list)
149             print("discard:", discard_list)           # print out discard
                                                    list
150         return data
151
152
153 # just for test
154 if __name__ == '__main__':
155     test_data = [
156         ['red', 25.6, 56, 1],
157         ['green', 33.3, 1, 1],
158         ['green', 2.5, 23, 0],
159         ['blue', 67.2, 111, 1],
160         ['red', 29.0, 34, 0],
161         ['yellow', 99.5, 78, 1],
162         ['yellow', 10.2, 23, 1],
163         ['yellow', 9.9, 30, 0],
164         ['blue', 67.0, 47, 0],
165         ['red', 41.8, 99, 1]
166     ]
167     test_attribute = ['color', 'average', 'age', 'class']
168     test_value_type = ['categorical', 'numerical', 'numerical', 'label']
169     test_data_after = pre_process(test_data, test_attribute,
                                    test_value_type)
170     print(test_data_after)

```

### 三、程序清单read.py

```

1  """
2  Description: Read initial dataset and decode it into a list. Here we
              replace all missing value and discretize
3              the numerical values.
4  Input: initial dataset stored in *.data file, and scheme description

```



```

        stored in *.names file.
5  Output: a data list after pre-processing.
6  Author: CBA Studio
7  """
8  import csv
9
10
11 # Read dataset and convert into a list.
12 # path: directory of *.data file.
13 def read_data(path):
14     data = []
15     with open(path, 'r') as csv_file:
16         reader = csv.reader(csv_file, delimiter=',')
17         for line in reader:
18             data.append(line)
19         while [] in data:
20             data.remove([])
21     return data
22
23
24 # Read scheme file *.names and write down attributes and value types.
25 # path: directory of *.names file.
26 def read_scheme(path):
27     with open(path, 'r') as csv_file:
28         reader = csv.reader(csv_file, delimiter=',')
29         attributes = next(reader)
30         value_type = next(reader)
31     return attributes, value_type
32
33
34 # convert string-type value into float-type.
35 # data: data list returned by read_data.
36 # value_type: list returned by read_scheme.
37 def str2numerical(data, value_type):
38     size = len(data)
39     columns = len(data[0])
40     for i in range(size):
41         for j in range(columns-1):
42             if value_type[j] == 'numerical' and data[i][j] != '?':
43                 data[i][j] = float(data[i][j])
44     return data
45
46
47 # Main method in this file, to get data list after processing and scheme
    list.
48 # data_path: tell where *.data file stores.
49 # scheme_path: tell where *.names file stores.
50 def read(data_path, scheme_path):

```

```

51     data = read_data(data_path)
52     attributes, value_type = read_scheme(scheme_path)
53     data = str2numerical(data, value_type)
54     return data, attributes, value_type
55
56
57 # just for test
58 if __name__ == '__main__':
59     import pre_processing
60
61     test_data_path = '/Users/liulizhi/Desktop/iris.data'
62     test_scheme_path = '/Users/liulizhi/Desktop/iris.names'
63     test_data, test_attributes, test_value_type = read(test_data_path,
64                                                         test_scheme_path)
65     result_data = pre_processing.pre_process(test_data, test_attributes,
66                                              test_value_type)
67     print(result_data)

```

#### 四、程序清单ruleitem.py

```

1  """
2  Description: Definition of class RuleItem, including condset, class label
3              (y in paper), condsupCount, rulesupCount,
4              support and confidence.
5  Input: condset which is a set of items, class label and the dataset.
6  Output: a ruleitem with its condsupCount, rulesupCount, support and
7          confidence.
8  Author: CBA Studio
9  Reference: https://www.cs.uic.edu/~hxiao/courses/cs594-slides.pdf
10 """
11
12 class RuleItem:
13     """
14     cond_set: a dict with following fashion:
15         {item name: value, item name: value, ...}
16         e.g.
17         {A: 1, B: 1} (A, B are name of columns, here called "item",
18                     and in our code should be numerical index
19                     but not string)
20     class_label: just to identify the class it belongs to.
21     dataset: a list returned by read method. (see read.py)
22     cond_sup_count, rule_sup_count, support and confidence are number.
23     """
24     def __init__(self, cond_set, class_label, dataset):
25         self.cond_set = cond_set
26         self.class_label = class_label
27         self.cond_sup_count, self.rule_sup_count = self._get_sup_count(
28             dataset)

```

```

26         self.support = self._get_support(len(dataset))
27         self.confidence = self._get_confidence()
28
29     # calculate condsupCount and rulesupCount
30     def _get_sup_count(self, dataset):
31         cond_sup_count = 0
32         rule_sup_count = 0
33         for case in dataset:
34             is_contained = True
35             for index in self.cond_set:
36                 if self.cond_set[index] != case[index]:
37                     is_contained = False
38                     break
39             if is_contained:
40                 cond_sup_count += 1
41                 if self.class_label == case[-1]:
42                     rule_sup_count += 1
43         return cond_sup_count, rule_sup_count
44
45     # calculate support count
46     def _get_support(self, dataset_size):
47         return self.rule_sup_count / dataset_size
48
49     # calculate confidence
50     def _get_confidence(self):
51         if self.cond_sup_count != 0:
52             return self.rule_sup_count / self.cond_sup_count
53         else:
54             return 0
55
56     # print out the ruleitem
57     def print(self):
58         cond_set_output = ''
59         for item in self.cond_set:
60             cond_set_output += '(' + str(item) + ', ' + str(self.cond_set
61                                     [item]) + '), '
62         cond_set_output = cond_set_output[:-2]
63         print('<({' + cond_set_output + '}', ' + str(self.cond_sup_count)
64               + '), (' +
65               '(class, ' + str(self.class_label) + '), ' + str(self.
66                   rule_sup_count) + ')>')
67
68     # print out rule
69     def print_rule(self):
70         cond_set_output = ''
71         for item in self.cond_set:
72             cond_set_output += '(' + str(item) + ', ' + str(self.cond_set
73                                     [item]) + '), '

```

```

70         cond_set_output = '{' + cond_set_output[:-2] + '}'
71         print(cond_set_output + ' -> (class, ' + str(self.class_label) +
              ')')
72
73
74 # just for test
75 if __name__ == '__main__':
76     cond_set = {0: 1, 1: 1}
77     class_label = 1
78     dataset = [[1, 1, 1], [1, 1, 1], [1, 2, 1], [2, 2, 1], [2, 2, 1],
79               [2, 2, 0], [2, 3, 0], [2, 3, 0], [1, 1, 0], [3, 2, 0]]
80     rule_item = RuleItem(cond_set, class_label, dataset)
81     rule_item.print()
82     rule_item.print_rule()
83     print('condsupCount =', rule_item.cond_sup_count)      # should be 3
84     print('rulesupCount =', rule_item.rule_sup_count)      # should be 2
85     print('support =', rule_item.support)                  # should be 0.2
86     print('confidence =', rule_item.confidence)            # should be 0.667

```

## 五、程序清单cba\_rg.py

```

1  """
2  Description: The implementation of CBA-RG algorithm, generating the
3               complete set of CARs (Class Association Rules).
4               We just follow up algorithm raised up in the paper without
5               improvement.
6  Input: a dataset got from pre_process (see pre_processing.py), minsup and
7         minconf
8  Output: CARs
9  Author: CBA Studio
10 Reference: https://www.cs.uic.edu/~hxiao/courses/cs594-slides.pdf
11 """
12 import ruleitem
13
14 class FrequentRuleitems:
15     """
16     A set of frequent k-ruleitems, just using set.
17     """
18     def __init__(self):
19         self.frequent_ruleitems_set = set()
20
21     # get size of set
22     def get_size(self):
23         return len(self.frequent_ruleitems_set)
24
25     # add a new ruleitem into set
26     def add(self, rule_item):
27         is_existed = False

```

```

26         for item in self.frequent_ruleitems_set:
27             if item.class_label == rule_item.class_label:
28                 if item.cond_set == rule_item.cond_set:
29                     is_existed = True
30                     break
31         if not is_existed:
32             self.frequent_ruleitems_set.add(rule_item)
33
34     # append set of ruleitems
35     def append(self, sets):
36         for item in sets.frequent_ruleitems:
37             self.add(item)
38
39     # print out all frequent ruleitems
40     def print(self):
41         for item in self.frequent_ruleitems_set:
42             item.print()
43
44
45 class Car:
46     """
47     Class Association Rules (Car). If some ruleitems has the same condset
48     , the ruleitem with the highest confidence is
49     chosen as the Possible Rule (PR). If there're more than one ruleitem
50     with the same highest confidence, we randomly
51     select one ruleitem.
52     """
53     def __init__(self):
54         self.rules = set()
55         self.pruned_rules = set()
56
57     # print out all rules
58     def print_rule(self):
59         for item in self.rules:
60             item.print_rule()
61
62     # print out all pruned rules
63     def print_pruned_rule(self):
64         for item in self.pruned_rules:
65             item.print_rule()
66
67     # add a new rule (frequent & accurate), save the ruleitem with the
68     # highest confidence when having the same condset
69     def _add(self, rule_item, minsup, minconf):
70         if rule_item.support >= minsup and rule_item.confidence >=
71             minconf:
72             if rule_item in self.rules:
73                 return

```

```

70         for item in self.rules:
71             if item.cond_set == rule_item.cond_set and item.
               confidence < rule_item.confidence:
72                 self.rules.remove(item)
73                 self.rules.add(rule_item)
74                 return
75             elif item.cond_set == rule_item.cond_set and item.
               confidence >= rule_item.confidence:
76                 return
77         self.rules.add(rule_item)
78
79     # convert frequent ruleitems into car
80     def gen_rules(self, frequent_ruleitems, minsup, minconf):
81         for item in frequent_ruleitems.frequent_ruleitems_set:
82             self._add(item, minsup, minconf)
83
84     # prune rules
85     def prune_rules(self, dataset):
86         for rule in self.rules:
87             pruned_rule = prune(rule, dataset)
88
89             is_existed = False
90             for rule in self.pruned_rules:
91                 if rule.class_label == pruned_rule.class_label:
92                     if rule.cond_set == pruned_rule.cond_set:
93                         is_existed = True
94                         break
95
96             if not is_existed:
97                 self.pruned_rules.add(pruned_rule)
98
99     # union new car into rules list
100    def append(self, car, minsup, minconf):
101        for item in car.rules:
102            self._add(item, minsup, minconf)
103
104
105    # try to prune rule
106    def prune(rule, dataset):
107        import sys
108        min_rule_error = sys.maxsize
109        pruned_rule = rule
110
111        # prune rule recursively
112        def find_prune_rule(this_rule):
113            nonlocal min_rule_error
114            nonlocal pruned_rule
115

```

```

116         # calculate how many errors the rule r make in the dataset
117     def errors_of_rule(r):
118         import cba_cb_m1
119
120         errors_number = 0
121         for case in dataset:
122             if cba_cb_m1.is_satisfy(case, r) == False:
123                 errors_number += 1
124         return errors_number
125
126     rule_error = errors_of_rule(this_rule)
127     if rule_error < min_rule_error:
128         min_rule_error = rule_error
129         pruned_rule = this_rule
130     this_rule_cond_set = list(this_rule.cond_set)
131     if len(this_rule_cond_set) >= 2:
132         for attribute in this_rule_cond_set:
133             temp_cond_set = dict(this_rule.cond_set)
134             temp_cond_set.pop(attribute)
135             temp_rule = ruleitem.RuleItem(temp_cond_set, this_rule.
136                                           class_label, dataset)
137             temp_rule_error = errors_of_rule(temp_rule)
138             if temp_rule_error <= min_rule_error:
139                 min_rule_error = temp_rule_error
140                 pruned_rule = temp_rule
141                 if len(temp_cond_set) >= 2:
142                     find_prune_rule(temp_rule)
143     find_prune_rule(rule)
144     return pruned_rule
145
146
147 # invoked by candidate_gen, join two items to generate candidate
148 def join(item1, item2, dataset):
149     if item1.class_label != item2.class_label:
150         return None
151     category1 = set(item1.cond_set)
152     category2 = set(item2.cond_set)
153     if category1 == category2:
154         return None
155     intersect = category1 & category2
156     for item in intersect:
157         if item1.cond_set[item] != item2.cond_set[item]:
158             return None
159     category = category1 | category2
160     new_cond_set = dict()
161     for item in category:
162         if item in category1:

```

```

163         new_cond_set[item] = item1.cond_set[item]
164     else:
165         new_cond_set[item] = item2.cond_set[item]
166     new_ruleitem = ruleitem.RuleItem(new_cond_set, item1.class_label,
167                                     dataset)
168     return new_ruleitem
169
170 # similar to Apriori-gen in algorithm Apriori
171 def candidate_gen(frequent_ruleitems, dataset):
172     returned_frequent_ruleitems = FrequentRuleitems()
173     for item1 in frequent_ruleitems.frequent_ruleitems_set:
174         for item2 in frequent_ruleitems.frequent_ruleitems_set:
175             new_ruleitem = join(item1, item2, dataset)
176             if new_ruleitem:
177                 returned_frequent_ruleitems.add(new_ruleitem)
178                 if returned_frequent_ruleitems.get_size() >= 1000:      #
179                     not allow to store more than 1000 ruleitems
180                     return returned_frequent_ruleitems
181     return returned_frequent_ruleitems
182
183 # main method, implementation of CBA-RG algorithm
184 def rule_generator(dataset, minsup, minconf):
185     frequent_ruleitems = FrequentRuleitems()
186     car = Car()
187
188     # get large 1-ruleitems and generate rules
189     class_label = set([x[-1] for x in dataset])
190     for column in range(0, len(dataset[0])-1):
191         distinct_value = set([x[column] for x in dataset])
192         for value in distinct_value:
193             cond_set = {column: value}
194             for classes in class_label:
195                 rule_item = ruleitem.RuleItem(cond_set, classes, dataset)
196                 if rule_item.support >= minsup:
197                     frequent_ruleitems.add(rule_item)
198     car.gen_rules(frequent_ruleitems, minsup, minconf)
199     cars = car
200
201     last_cars_number = 0
202     current_cars_number = len(cars.rules)
203     while frequent_ruleitems.get_size() > 0 and current_cars_number <=
204         2000 and \
205         (current_cars_number - last_cars_number) >= 10:
206         candidate = candidate_gen(frequent_ruleitems, dataset)
207         frequent_ruleitems = FrequentRuleitems()
208         car = Car()

```



```

208         for item in candidate.frequent_ruleitems_set:
209             if item.support >= minsup:
210                 frequent_ruleitems.add(item)
211             car.gen_rules(frequent_ruleitems, minsup, minconf)
212             cars.append(car, minsup, minconf)
213             last_cars_number = current_cars_number
214             current_cars_number = len(cars.rules)
215
216     return cars
217
218
219 # just for test
220 if __name__ == "__main__":
221     dataset = [[1, 1, 1], [1, 1, 1], [1, 2, 1], [2, 2, 1], [2, 2, 1],
222               [2, 2, 0], [2, 3, 0], [2, 3, 0], [1, 1, 0], [3, 2, 0]]
223     minsup = 0.15
224     minconf = 0.6
225     cars = rule_generator(dataset, minsup, minconf)
226
227     print("CARs:")
228     cars.print_rule()
229
230     print("prCARs:")
231     cars.prune_rules(dataset)
232     cars.print_pruned_rule()

```

## 六、程序清单cba\_cb\_m1.py

```

1  """
2  Description: The implementation of a naive algorithm for CBA-CB: M1. The
3              class Classifier includes the set of selected
4              rules and default_class, which can be expressed as <r1, r2, ..., rn,
5              default_class>. Method classifier_builder_m1
6              is the main method to implement CBA-CB: M1.
7  Input: a set of CARs generated from rule_generator (see cab_rg.py) and a
8         dataset got from pre_process
9         (see pre_processing.py)
10 Output: a classifier
11 Author: CBA Studio
12 Reference: http://www.docin.com/p-586554186.html
13 """
14
15 import cba_rg
16 from functools import cmp_to_key
17 import sys
18
19
20 # check the rule whether covers the data case (a line in table with the
21 # class label at the end of list) or not
22
23 # if covers (LHS of rule is the same as the data case, and they belongs

```

```

        to the same class), return True;
18 # else if LHSs are the same while the class labels are different, return
    False;
19 # else (LHSs are different), return None
20 def is_satisfy(datacase, rule):
21     for item in rule.cond_set:
22         if datacase[item] != rule.cond_set[item]:
23             return None
24     if datacase[-1] == rule.class_label:
25         return True
26     else:
27         return False
28
29
30 class Classifier:
31     """
32     This class is our classifier. The rule_list and default_class are
        useful for outer code.
33     """
34     def __init__(self):
35         self.rule_list = list()
36         self.default_class = None
37         self._error_list = list()
38         self._default_class_list = list()
39
40         # insert a rule into rule_list, then choose a default class, and
        calculate the errors (see line 8, 10 & 11)
41     def insert(self, rule, dataset):
42         self.rule_list.append(rule)                # insert r at the end of
        C
43         self._select_default_class(dataset)         # select a default class
        for the current C
44         self._compute_error(dataset)                # compute the total
        number of errors of C
45
46         # select the majority class in the remaining data
47     def _select_default_class(self, dataset):
48         class_column = [x[-1] for x in dataset]
49         class_label = set(class_column)
50         max = 0
51         current_default_class = None
52         for label in class_label:
53             if class_column.count(label) >= max:
54                 max = class_column.count(label)
55                 current_default_class = label
56         self._default_class_list.append(current_default_class)
57
58         # compute the sum of errors

```

```

59     def _compute_error(self, dataset):
60         if len(dataset) <= 0:
61             self._error_list.append(sys.maxsize)
62             return
63
64         error_number = 0
65
66         # the number of errors that have been made by all the selected
67         # rules in C
68         for case in dataset:
69             is_cover = False
70             for rule in self.rule_list:
71                 if is_satisfy(case, rule):
72                     is_cover = True
73                     break
74             if not is_cover:
75                 error_number += 1
76
77         # the number of errors to be made by the default class in the
78         # training set
79         class_column = [x[-1] for x in dataset]
80         error_number += len(class_column) - class_column.count(self.
81             _default_class_list[-1])
82         self._error_list.append(error_number)
83
84     # see line 14 and 15, to get the final classifier
85     def discard(self):
86         # find the first rule p in C with the lowest total number of
87         # errors and drop all the rules after p in C
88         index = self._error_list.index(min(self._error_list))
89         self.rule_list = self.rule_list[: (index+1)]
90         self._error_list = None
91
92         # assign the default class associated with p to default_class
93         self.default_class = self._default_class_list[index]
94         self._default_class_list = None
95
96     # just print out all selected rules and default class in our
97     # classifier
98     def print(self):
99         for rule in self.rule_list:
100             rule.print_rule()
101         print("default_class:", self.default_class)
102
103     # sort the set of generated rules car according to the relation ">",
104     # return the sorted rule list
105     def sort(car):

```

```

101     def cmp_method(a, b):
102         if a.confidence < b.confidence:      # 1. the confidence of ri >
            rj
103             return 1
104         elif a.confidence == b.confidence:
105             if a.support < b.support:        # 2. their confidences are
                the same, but support of ri > rj
106                 return 1
107             elif a.support == b.support:
108                 if len(a.cond_set) < len(b.cond_set):    # 3. both
                    confidence & support are the same, ri earlier than rj
109                     return -1
110                 elif len(a.cond_set) == len(b.cond_set):
111                     return 0
112                 else:
113                     return 1
114             else:
115                 return -1
116         else:
117             return -1
118
119     rule_list = list(car.rules)
120     rule_list.sort(key=cmp_to_key(cmp_method))
121     return rule_list
122
123
124 # main method of CBA-CB: M1
125 def classifier_builder_m1(cars, dataset):
126     classifier = Classifier()
127     cars_list = sort(cars)
128     for rule in cars_list:
129         temp = []
130         mark = False
131         for i in range(len(dataset)):
132             is_satisfy_value = is_satisfy(dataset[i], rule)
133             if is_satisfy_value is not None:
134                 temp.append(i)
135                 if is_satisfy_value:
136                     mark = True
137         if mark:
138             temp_dataset = list(dataset)
139             for index in temp:
140                 temp_dataset[index] = []
141             while [] in temp_dataset:
142                 temp_dataset.remove([])
143             dataset = temp_dataset
144             classifier.insert(rule, dataset)
145     classifier.discard()

```

```

146     return classifier
147
148
149 # just for test
150 if __name__ == '__main__':
151     dataset = [[1, 1, 1], [1, 1, 1], [1, 2, 1], [2, 2, 1], [2, 2, 1],
152               [2, 2, 0], [2, 3, 0], [2, 3, 0], [1, 1, 0], [3, 2, 0]]
153     minsup = 0.15
154     minconf = 0.6
155     cars = cba_rg.rule_generator(dataset, minsup, minconf)
156     classifier = classifier_builder_m1(cars, dataset)
157     classifier.print()
158
159     print()
160     dataset = [[1, 1, 1], [1, 1, 1], [1, 2, 1], [2, 2, 1], [2, 2, 1],
161               [2, 2, 0], [2, 3, 0], [2, 3, 0], [1, 1, 0], [3, 2, 0]]
162     cars.prune_rules(dataset)
163     cars.rules = cars.pruned_rules
164     classifier = classifier_builder_m1(cars, dataset)
165     classifier.print()

```

## 七、程序清单 cba\_cb\_m2.py

```

1  """
2  Description: The following code implements an improved version of the
3               algorithm, called CBA-CB: M2. It contains three
4               stages. For stage 1, we scan the whole database, to find the cRule
5               and wRule, get the set Q, U and A at the same
6               time. In stage 2, for each case d that we could not decide which rule
7               should cover it in stage 1, we go through d
8               again to find all rules that classify it wrongly and have a higher
9               precedence than the corresponding cRule of d.
10              Finally, in stage 3, we choose the final set of rules to form our
11              final classifier.
12  Input: a set of CARs generated from rule_generator (see cab_rg.py) and a
13         dataset got from pre_process
14         (see pre_processing.py)
15  Output: a classifier
16  Author: CBA Studio
17  """
18
19  import ruleitem
20  import cba_cb_m1
21  from functools import cmp_to_key
22
23
24  class Classifier_m2:
25      """
26      The definition of classifier formed in CBA-CB: M2. It contains a list
27      of rules order by their precedence, a default

```

```

20     class label. The other member are private and useless for outer code.
21     """
22     def __init__(self):
23         self.rule_list = list()
24         self.default_class = None
25         self._default_class_list = list()
26         self._total_errors_list = list()
27
28     # insert a new rule into classifier
29     def add(self, rule, default_class, total_errors):
30         self.rule_list.append(rule)
31         self._default_class_list.append(default_class)
32         self._total_errors_list.append(total_errors)
33
34     # discard those rules that introduce more errors. See line 18-20, CBA
35     # -CB: M2 (Stage 3).
36     def discard(self):
37         index = self._total_errors_list.index(min(self._total_errors_list
38             ))
39         self.rule_list = self.rule_list[: (index + 1)]
40         self._total_errors_list = None
41
42         self.default_class = self._default_class_list[index]
43         self._default_class_list = None
44
45     # just print out rules and default class label
46     def print(self):
47         for rule in self.rule_list:
48             rule.print_rule()
49         print("default_class:", self.default_class)
50
51 class Rule(ruleitem.RuleItem):
52     """
53     A class inherited from RuleItem, adding classCasesCovered and replace
54     field.
55     """
56     def __init__(self, cond_set, class_label, dataset):
57         ruleitem.RuleItem.__init__(self, cond_set, class_label, dataset)
58         self._init_classCasesCovered(dataset)
59         self.replace = set()
60
61     # initialize the classCasesCovered field
62     def _init_classCasesCovered(self, dataset):
63         class_column = [x[-1] for x in dataset]
64         class_label = set(class_column)
65         self.classCasesCovered = dict((x, 0) for x in class_label)

```

```

65
66 # convert ruleitem of class RuleItem to rule of class Rule
67 def ruleitem2rule(rule_item, dataset):
68     rule = Rule(rule_item.cond_set, rule_item.class_label, dataset)
69     return rule
70
71
72 # finds the highest precedence rule that covers the data case d from the
    set of rules having the same class as d.
73 def maxCoverRule_correct(cars_list, data_case):
74     for i in range(len(cars_list)):
75         if cars_list[i].class_label == data_case[-1]:
76             if cba_cb_m1.is_satisfy(data_case, cars_list[i]):
77                 return i
78     return None
79
80
81 # finds the highest precedence rule that covers the data case d from the
    set of rules having the different class as d.
82 def maxCoverRule_wrong(cars_list, data_case):
83     for i in range(len(cars_list)):
84         if cars_list[i].class_label != data_case[-1]:
85             temp_data_case = data_case[:-1]
86             temp_data_case.append(cars_list[i].class_label)
87             if cba_cb_m1.is_satisfy(temp_data_case, cars_list[i]):
88                 return i
89     return None
90
91
92 # compare two rule, return the precedence.
93 # -1: rule1 < rule2, 0: rule1 < rule2 (randomly here), 1: rule1 > rule2
94 def compare(rule1, rule2):
95     if rule1 is None and rule2 is not None:
96         return -1
97     elif rule1 is None and rule2 is None:
98         return 0
99     elif rule1 is not None and rule2 is None:
100         return 1
101
102     if rule1.confidence < rule2.confidence:      # 1. the confidence of ri
        > rj
103         return -1
104     elif rule1.confidence == rule2.confidence:
105         if rule1.support < rule2.support:      # 2. their confidences
            are the same, but support of ri > rj
106             return -1
107         elif rule1.support == rule2.support:
108             if len(rule1.cond_set) < len(rule2.cond_set):    # 3.

```

```

        confidence & support are the same, ri earlier than rj
109         return 1
110     elif len(rule1.cond_set) == len(rule2.cond_set):
111         return 0
112     else:
113         return -1
114     else:
115         return 1
116 else:
117     return 1
118
119
120 # finds all the rules in u that wrongly classify the data case and have
    higher precedences than that of its cRule.
121 def allCoverRules(u, data_case, c_rule, cars_list):
122     w_set = set()
123     for rule_index in u:
124         # have higher precedences than cRule
125         if compare(cars_list[rule_index], c_rule) > 0:
126             # wrongly classify the data case
127             if cba_cb_m1.is_satisfy(data_case, cars_list[rule_index]) ==
                False:
128                 w_set.add(rule_index)
129     return w_set
130
131
132 # counts the number of training cases in each class
133 def compClassDistr(dataset):
134     class_distr = dict()
135
136     if len(dataset) <= 0:
137         class_distr = None
138
139     dataset_without_null = dataset
140     while [] in dataset_without_null:
141         dataset_without_null.remove([])
142
143     class_column = [x[-1] for x in dataset_without_null]
144     class_label = set(class_column)
145     for c in class_label:
146         class_distr[c] = class_column.count(c)
147     return class_distr
148
149
150 # sort the rule list order by precedence
151 def sort_with_index(q, cars_list):
152     def cmp_method(a, b):
153         # 1. the confidence of ri > rj

```



```

154         if cars_list[a].confidence < cars_list[b].confidence:
155             return 1
156         elif cars_list[a].confidence == cars_list[b].confidence:
157             # 2. their confidences are the same, but support of ri > rj
158             if cars_list[a].support < cars_list[b].support:
159                 return 1
160             elif cars_list[a].support == cars_list[b].support:
161                 # 3. both confidence & support are the same, ri earlier
162                 # than rj
163                 if len(cars_list[a].cond_set) < len(cars_list[b].cond_set):
164                     return -1
165                 elif len(cars_list[a].cond_set) == len(cars_list[b].
166                     cond_set):
167                     return 0
168                 else:
169                     return 1
170             else:
171                 return -1
172
173     rule_list = list(q)
174     rule_list.sort(key=cmp_to_key(cmp_method))
175     return set(rule_list)
176
177
178 # get how many errors the rule wrongly classify the data case
179 def errorsOfRule(rule, dataset):
180     error_number = 0
181     for case in dataset:
182         if case:
183             if cba_cb_m1.is_satisfy(case, rule) == False:
184                 error_number += 1
185     return error_number
186
187
188 # choose the default class (majority class in remaining dataset)
189 def selectDefault(class_distribution):
190     if class_distribution is None:
191         return None
192
193     max = 0
194     default_class = None
195     for index in class_distribution:
196         if class_distribution[index] > max:
197             max = class_distribution[index]
198             default_class = index

```

```

199     return default_class
200
201
202 # count the number of errors that the default class will make in the
    remaining training data
203 def defErr(default_class, class_distribution):
204     if class_distribution is None:
205         import sys
206         return sys.maxsize
207
208     error = 0
209     for index in class_distribution:
210         if index != default_class:
211             error += class_distribution[index]
212     return error
213
214
215 # main method, implement the whole classifier builder
216 def classifier_builder_m2(cars, dataset):
217     classifier = Classifier_m2()
218
219     cars_list = cba_cb_m1.sort(cars)
220     for i in range(len(cars_list)):
221         cars_list[i] = ruleitem2rule(cars_list[i], dataset)
222
223     # stage 1
224     q = set()
225     u = set()
226     a = set()
227     mark_set = set()
228     for i in range(len(dataset)):
229         c_rule_index = maxCoverRule_correct(cars_list, dataset[i])
230         w_rule_index = maxCoverRule_wrong(cars_list, dataset[i])
231         if c_rule_index is not None:
232             u.add(c_rule_index)
233         if c_rule_index:
234             cars_list[c_rule_index].classCasesCovered[dataset[i][-1]] +=
                1
235         if c_rule_index and w_rule_index:
236             if compare(cars_list[c_rule_index], cars_list[w_rule_index])
                > 0:
237                 q.add(c_rule_index)
238                 mark_set.add(c_rule_index)
239             else:
240                 a.add((i, dataset[i][-1], c_rule_index, w_rule_index))
241         elif c_rule_index is None and w_rule_index is not None:
242             a.add((i, dataset[i][-1], c_rule_index, w_rule_index))
243

```

```

244     # stage 2
245     for entry in a:
246         if cars_list[entry[3]] in mark_set:
247             if entry[2] is not None:
248                 cars_list[entry[2]].classCasesCovered[entry[1]] -= 1
249                 cars_list[entry[3]].classCasesCovered[entry[1]] += 1
250             else:
251                 if entry[2] is not None:
252                     w_set = allCoverRules(u, dataset[entry[0]], cars_list[
253                         entry[2]], cars_list)
254                 else:
255                     w_set = allCoverRules(u, dataset[entry[0]], None,
256                         cars_list)
257                 for w in w_set:
258                     cars_list[w].replace.add((entry[2], entry[0], entry[1]))
259                     cars_list[w].classCasesCovered[entry[1]] += 1
260                 q |= w_set
261
262     # stage 3
263     rule_errors = 0
264     q = sort_with_index(q, cars_list)
265     data_cases_covered = list([False] * len(dataset))
266     for r_index in q:
267         if cars_list[r_index].classCasesCovered[cars_list[r_index].
268             class_label] != 0:
269             for entry in cars_list[r_index].replace:
270                 if data_cases_covered[entry[1]]:
271                     cars_list[r_index].classCasesCovered[entry[2]] -= 1
272                 else:
273                     if entry[0] is not None:
274                         cars_list[entry[0]].classCasesCovered[entry[2]]
275                             -= 1
276             for i in range(len(dataset)):
277                 datacase = dataset[i]
278                 if datacase:
279                     is_satisfy_value = cba_cb_m1.is_satisfy(datacase,
280                         cars_list[r_index])
281                     if is_satisfy_value:
282                         dataset[i] = []
283                         data_cases_covered[i] = True
284             rule_errors += errorsOfRule(cars_list[r_index], dataset)
285             class_distribution = compClassDistr(dataset)
286             default_class = selectDefault(class_distribution)
287             default_errors = defErr(default_class, class_distribution)
288             total_errors = rule_errors + default_errors
289             classifier.add(cars_list[r_index], default_class,
290                 total_errors)
291     classifier.discard()

```

```

286
287     return classifier
288
289
290 # just for test
291 if __name__ == "__main__":
292     import cba_rg
293
294     dataset = [[1, 1, 1], [1, 1, 1], [1, 2, 1], [2, 2, 1], [2, 2, 1],
295                [2, 2, 0], [2, 3, 0], [2, 3, 0], [1, 1, 0], [3, 2, 0]]
296     minsup = 0.15
297     minconf = 0.6
298     cars = cba_rg.rule_generator(dataset, minsup, minconf)
299     classifier = classifier_builder_m2(cars, dataset)
300     classifier.print()
301
302     print()
303     dataset = [[1, 1, 1], [1, 1, 1], [1, 2, 1], [2, 2, 1], [2, 2, 1],
304                [2, 2, 0], [2, 3, 0], [2, 3, 0], [1, 1, 0], [3, 2, 0]]
305     cars.prune_rules(dataset)
306     cars.rules = cars.pruned_rules
307     classifier = classifier_builder_m2(cars, dataset)
308     classifier.print()

```

## 八、程序清单validation.py

```

1  """
2  Description: This is our experimental code. We provide 4 test modes for
3               experiments:
4               10-fold cross-validations on CBA (M1) without pruning
5               10-fold cross-validations on CBA (M1) with pruning
6               10-fold cross-validations on CBA (M2) without pruning
7               10-fold cross-validations on CBA (M2) with pruning
8  Input: the relative directory path of data file and scheme file
9  Output: the experimental results (similar to Table 1: Experiment Results
10         in this paper)
11 Author: CBA Studio
12 """
13 from read import read
14 from pre_processing import pre_process
15 from cba_rg import rule_generator
16 from cba_cb_m1 import classifier_builder_m1
17 from cba_cb_m1 import is_satisfy
18 from cba_cb_m2 import classifier_builder_m2
19 import time
20 import random
21 # calculate the error rate of the classifier on the dataset

```

```

22 def get_error_rate(classifier, dataset):
23     size = len(dataset)
24     error_number = 0
25     for case in dataset:
26         is_satisfy_value = False
27         for rule in classifier.rule_list:
28             is_satisfy_value = is_satisfy(case, rule)
29             if is_satisfy_value == True:
30                 break
31         if is_satisfy_value == False:
32             if classifier.default_class != case[-1]:
33                 error_number += 1
34     return error_number / size
35
36
37 # 10-fold cross-validations on CBA (M1) without pruning
38 def cross_validate_ml_without_prune(data_path, scheme_path, minsup=0.01,
39     minconf=0.5):
40     data, attributes, value_type = read(data_path, scheme_path)
41     random.shuffle(data)
42     dataset = pre_process(data, attributes, value_type)
43
44     block_size = int(len(dataset) / 10)
45     split_point = [k * block_size for k in range(0, 10)]
46     split_point.append(len(dataset))
47
48     cba_rg_total_runtime = 0
49     cba_cb_total_runtime = 0
50     total_car_number = 0
51     total_classifier_rule_num = 0
52     error_total_rate = 0
53
54     for k in range(len(split_point)-1):
55         print("\nRound %d:" % k)
56
57         training_dataset = dataset[:split_point[k]] + dataset[split_point
58             [k+1]:]
59         test_dataset = dataset[split_point[k]:split_point[k+1]]
60
61         start_time = time.time()
62         cars = rule_generator(training_dataset, minsup, minconf)
63         end_time = time.time()
64         cba_rg_runtime = end_time - start_time
65         cba_rg_total_runtime += cba_rg_runtime
66
67         start_time = time.time()
68         classifier_ml = classifier_builder_ml(cars, training_dataset)
69         end_time = time.time()

```

```

68     cba_cb_runtime = end_time - start_time
69     cba_cb_total_runtime += cba_cb_runtime
70
71     error_rate = get_error_rate(classifier_m1, test_dataset)
72     error_total_rate += error_rate
73
74     total_car_number += len(cars.rules)
75     total_classifier_rule_num += len(classifier_m1.rule_list)
76
77     print("CBA's error rate without pruning: %.11f%%" % (error_rate *
78         100))
79     print("No. of CARs without pruning: %d" % len(cars.rules))
80     print("CBA-RG's run time without pruning: %.21f s" %
81         cba_rg_runtime)
82     print("CBA-CB M1's run time without pruning: %.21f s" %
83         cba_cb_runtime)
84     print("No. of rules in classifier of CBA-CB M1 without pruning: %
85         d" % len(classifier_m1.rule_list))
86
87     print("\nAverage CBA's error rate without pruning: %.11f%%" % (
88         error_total_rate / 10 * 100))
89     print("Average No. of CARs without pruning: %d" % int(
90         total_car_number / 10))
91     print("Average CBA-RG's run time without pruning: %.21f s" % (
92         cba_rg_total_runtime / 10))
93     print("Average CBA-CB M1's run time without pruning: %.21f s" % (
94         cba_cb_total_runtime / 10))
95     print("Average No. of rules in classifier of CBA-CB M1 without
96         pruning: %d" % int(total_classifier_rule_num / 10))
97
98
99
100 # 10-fold cross-validations on CBA (M1) with pruning
101 def cross_validate_m1_with_prune(data_path, scheme_path, minsup=0.01,
102     minconf=0.5):
103     data, attributes, value_type = read(data_path, scheme_path)
104     random.shuffle(data)
105     dataset = pre_process(data, attributes, value_type)
106
107     block_size = int(len(dataset) / 10)
108     split_point = [k * block_size for k in range(0, 10)]
109     split_point.append(len(dataset))
110
111     cba_rg_total_runtime = 0
112     cba_cb_total_runtime = 0
113     total_car_number = 0
114     total_classifier_rule_num = 0
115     error_total_rate = 0

```

```

106     for k in range(len(split_point)-1):
107         print("\nRound %d:" % k)
108
109         training_dataset = dataset[:split_point[k]] + dataset[split_point
110             [k+1]:]
111         test_dataset = dataset[split_point[k]:split_point[k+1]]
112
113         start_time = time.time()
114         cars = rule_generator(training_dataset, minsup, minconf)
115         cars.prune_rules(training_dataset)
116         cars.rules = cars.pruned_rules
117         end_time = time.time()
118         cba_rg_runtime = end_time - start_time
119         cba_rg_total_runtime += cba_rg_runtime
120
121         start_time = time.time()
122         classifier_m1 = classifier_builder_m1(cars, training_dataset)
123         end_time = time.time()
124         cba_cb_runtime = end_time - start_time
125         cba_cb_total_runtime += cba_cb_runtime
126
127         error_rate = get_error_rate(classifier_m1, test_dataset)
128         error_total_rate += error_rate
129
130         total_car_number += len(cars.rules)
131         total_classifier_rule_num += len(classifier_m1.rule_list)
132
133         print("CBA's error rate with pruning: %.11f%%" % (error_rate *
134             100))
135         print("No. of CARs with pruning: %d" % len(cars.rules))
136         print("CBA-RG's run time with pruning: %.21f s" % cba_rg_runtime)
137         print("CBA-CB M1's run time with pruning: %.21f s" %
138             cba_cb_runtime)
139         print("No. of rules in classifier of CBA-CB M1 with pruning: %d"
140             % len(classifier_m1.rule_list))
141
142     print("\nAverage CBA's error rate with pruning: %.11f%%" % (
143         error_total_rate / 10 * 100))
144     print("Average No. of CARs with pruning: %d" % int(total_car_number /
145         10))
146     print("Average CBA-RG's run time with pruning: %.21f s" % (
147         cba_rg_total_runtime / 10))
148     print("Average CBA-CB M1's run time with pruning: %.21f s" % (
149         cba_cb_total_runtime / 10))
150     print("Average No. of rules in classifier of CBA-CB M1 with pruning:
151         %d" % int(total_classifier_rule_num / 10))

```

```

145 # 10-fold cross-validations on CBA (M2) without pruning
146 def cross_validate_m2_without_prune(data_path, scheme_path, minsup=0.01,
    minconf=0.5):
147     data, attributes, value_type = read(data_path, scheme_path)
148     random.shuffle(data)
149     dataset = pre_process(data, attributes, value_type)
150
151     block_size = int(len(dataset) / 10)
152     split_point = [k * block_size for k in range(0, 10)]
153     split_point.append(len(dataset))
154
155     cba_rg_total_runtime = 0
156     cba_cb_total_runtime = 0
157     total_car_number = 0
158     total_classifier_rule_num = 0
159     error_total_rate = 0
160
161     for k in range(len(split_point)-1):
162         print("\nRound %d:" % k)
163
164         training_dataset = dataset[:split_point[k]] + dataset[split_point
            [k+1]:]
165         test_dataset = dataset[split_point[k]:split_point[k+1]]
166
167         start_time = time.time()
168         cars = rule_generator(training_dataset, minsup, minconf)
169         end_time = time.time()
170         cba_rg_runtime = end_time - start_time
171         cba_rg_total_runtime += cba_rg_runtime
172
173         start_time = time.time()
174         classifier_m2 = classifier_builder_m2(cars, training_dataset)
175         end_time = time.time()
176         cba_cb_runtime = end_time - start_time
177         cba_cb_total_runtime += cba_cb_runtime
178
179         error_rate = get_error_rate(classifier_m2, test_dataset)
180         error_total_rate += error_rate
181
182         total_car_number += len(cars.rules)
183         total_classifier_rule_num += len(classifier_m2.rule_list)
184
185         print("CBA's error rate without pruning: %.11f%%" % (error_rate *
            100))
186         print("No. of CARs without pruning: %d" % len(cars.rules))
187         print("CBA-RG's run time without pruning: %.21f s" %
            cba_rg_runtime)
188         print("CBA-CB M2's run time without pruning: %.21f s" %

```



```

        cba_cb_runtime)
189     print("No. of rules in classifier of CBA-CB M2 without pruning: %
        d" % len(classifier_m2.rule_list))
190
191     print("\nAverage CBA's error rate without pruning: %.11f%%" % (
        error_total_rate / 10 * 100))
192     print("Average No. of CARs without pruning: %d" % int(
        total_car_number / 10))
193     print("Average CBA-RG's run time without pruning: %.21f s" % (
        cba_rg_total_runtime / 10))
194     print("Average CBA-CB M2's run time without pruning: %.21f s" % (
        cba_cb_total_runtime / 10))
195     print("Average No. of rules in classifier of CBA-CB M2 without
        pruning: %d" % int(total_classifier_rule_num / 10))
196
197
198 # 10-fold cross-validations on CBA (M2) with pruning
199 def cross_validate_m2_with_prune(data_path, scheme_path, minsup=0.01,
        minconf=0.5):
200     data, attributes, value_type = read(data_path, scheme_path)
201     random.shuffle(data)
202     dataset = pre_process(data, attributes, value_type)
203
204     block_size = int(len(dataset) / 10)
205     split_point = [k * block_size for k in range(0, 10)]
206     split_point.append(len(dataset))
207
208     cba_rg_total_runtime = 0
209     cba_cb_total_runtime = 0
210     total_car_number = 0
211     total_classifier_rule_num = 0
212     error_total_rate = 0
213
214     for k in range(len(split_point)-1):
215         print("\nRound %d:" % k)
216
217         training_dataset = dataset[:split_point[k]] + dataset[split_point
            [k+1]:]
218         test_dataset = dataset[split_point[k]:split_point[k+1]]
219
220         start_time = time.time()
221         cars = rule_generator(training_dataset, minsup, minconf)
222         cars.prune_rules(training_dataset)
223         cars.rules = cars.pruned_rules
224         end_time = time.time()
225         cba_rg_runtime = end_time - start_time
226         cba_rg_total_runtime += cba_rg_runtime
227

```

```

228     start_time = time.time()
229     classifier_m2 = classifier_builder_m2(cars, training_dataset)
230     end_time = time.time()
231     cba_cb_runtime = end_time - start_time
232     cba_cb_total_runtime += cba_cb_runtime
233
234     error_rate = get_error_rate(classifier_m2, test_dataset)
235     error_total_rate += error_rate
236
237     total_car_number += len(cars.rules)
238     total_classifier_rule_num += len(classifier_m2.rule_list)
239
240     print("CBA's error rate with pruning: %.11f%%" % (error_rate *
        100))
241     print("No. of CARs without pruning: %d" % len(cars.rules))
242     print("CBA-RG's run time with pruning: %.21f s" % cba_rg_runtime)
243     print("CBA-CB M2's run time with pruning: %.21f s" %
        cba_cb_runtime)
244     print("No. of rules in classifier of CBA-CB M2 with pruning: %d"
        % len(classifier_m2.rule_list))
245
246     print("\nAverage CBA's error rate with pruning: %.11f%%" % (
        error_total_rate / 10 * 100))
247     print("Average No. of CARs with pruning: %d" % int(total_car_number /
        10))
248     print("Average CBA-RG's run time with pruning: %.21f s" % (
        cba_rg_total_runtime / 10))
249     print("Average CBA-CB M2's run time with pruning: %.21f s" % (
        cba_cb_total_runtime / 10))
250     print("Average No. of rules in classifier of CBA-CB M2 with pruning:
        %d" % int(total_classifier_rule_num / 10))
251
252
253 # test entry goes here
254 if __name__ == "__main__":
255     # using the relative path, all data sets are stored in datasets
        directory
256     test_data_path = 'datasets/australian.data'
257     test_scheme_path = 'datasets/australian.names'
258
259     # just choose one mode to experiment by removing one line comment and
        running
260     cross_validate_m1_without_prune(test_data_path, test_scheme_path)
261     # cross_validate_m1_with_prune(test_data_path, test_scheme_path)
262     # cross_validate_m2_without_prune(test_data_path, test_scheme_path)
263     # cross_validate_m1_with_prune(test_data_path, test_scheme_path)

```

## 附录 B 本文使用的数据集格式与示例

本文中的数据集全部选取自 UCI 机器学习资料库 [2]，但是由于网站中数据集格式千差万别，不便于我们的程序读取和处理数据，所以我们在不改变数据内容的前提下，对所采用的 30 个数据集进行格式归一化处理。

对于一个数据集，我们需要两份文件。一份以 `.data` 为后缀，里面存储的是以逗号分隔的 `csv` 格式的样本观测数据，每一行为一个样本，每一列对应一个属性，最后一列为该样本所属的类的标号。若某一样本的某一属性值缺失，则将其用英文问号“?”标注。另一份以 `.names` 为后缀，描述了该数据集的格式。第一行为属性列的名称，以逗号分隔，最后一列命名为“`class`”，以表示所属类。第二行为各属性值的性质，若属性值为连续型（又称“数值型”），则标记为“`numerical`”；若属性值为离散型（又称“标签型”），则标记为“`categorical`”；最后一列为类标号，记为“`label`”。必须保证这两行的列数相等，并且和 `.data` 文件的列数一致。

以著名的 `iris` 数据集为例。下面给出的是 `iris.data` 文件中前 10 行数据，它们对应 10 个样本，前 4 列为观测到的属性值，最后一列对应这个样本所属的类。

### 文件清单 1 `iris.data` 文件的前 10 行

```
5.1,3.5,1.4,0.2,Iris-setosa
4.9,3.0,1.4,0.2,Iris-setosa
4.7,3.2,1.3,0.2,Iris-setosa
4.6,3.1,1.5,0.2,Iris-setosa
5.0,3.6,1.4,0.2,Iris-setosa
5.4,3.9,1.7,0.4,Iris-setosa
4.6,3.4,1.4,0.3,Iris-setosa
5.0,3.4,1.5,0.2,Iris-setosa
4.4,2.9,1.4,0.2,Iris-setosa
4.9,3.1,1.5,0.1,Iris-setosa
```

接着给出的是 `iris.names` 文件的内容。第 1 行前 4 列均为各属性的名称；最后 1 列为“`class`”，表明是类标号。第 2 行前四列全为“`numerical`”，表示这些属性值都是连续型；最后 1 列的“`label`”表示其为类的标号。

### 文件清单 2 `iris.names` 文件

```
sepal length,sepal width,petal length,petal width,class
numerical,numerical,numerical,numerical,label
```

其他的数据集均按照上述规则进行归一化处理。运用这些数据集，可以对我们的程序进行实例验证，以测试我们实现的性能。