

# 基于改良蚁群算法的神经网络分类规则提取<sup>①</sup>

许海波, 刘端阳, 胡同森

(浙江工业大学 计算机科学与技术学院, 杭州 310023)

**摘要:** 在数据挖掘领域, 分类获得了很大的关注度, 其主要目的是预测数据对象的所属类别。分类方法可分为基于规则和不基于规则两大类, 其中神经网络由于在预测、从经验中学习、从先前样本中泛化等方面的优秀表现, 使其成为分类领域的一个重要的方法, 并往往能够获得很高的分类准确性, 然而其非常有限的解释能力成为了制约其应用的一大缺陷。提出了一种基于改良蚁群算法的神经网络分类规则提取方法, 通过改良的蚁群算法来填补神经网络有限的解释能力, 从数据中提取出分类规则。实验证明, 该方法能够很好的辅助神经网络, 从要分类的数据中获取规则。

**关键词:** 数据分类; 数据挖掘; 规则提取; 蚁群算法; 神经网络

## Rules Extraction from Artificial Neural Networks for Classification Based Improved Ant Colony Algorithm

XU Hai-Bo, LIU Duan-Yang, HU Tong-Sen

(College of Computer Science and Technology, Zhejiang University of Technology, Hangzhou 310023, China)

**Abstract:** Classification obtains great concern in the field of data mining. Its main purpose is to predict the classification of data objects. Classification can be divided into two major categories of rule-based and non-rule-based, however because of the excellent performance that artificial neural network(ANN) can obtain from prediction, studying from experience and generalizing from the previous samples, making it an important method of classification. Although ANNs can achieve high classification accuracy, their explanation capability is very limited, as to restrict its application. This paper presents an improved ant colony algorithm based on ANNs classification rule extraction method, an improved ant colony algorithm is to help solve the ANN's limited explanation capability to extract rules from the data. Experiments show that this approach could coordinate neural network to obtain rules of classified data well.

**Key words:** data classification; data mining; rules extraction; ant colony algorithm; artificial neural networks

数据挖掘即分析所获得的数据集来找出不为人所关注的联系和以对数据所有者有用且易于理解的新颖的途径来概括数据, 在经济, 科学等大量领域获得了广泛应用。它能够帮助用户仅仅需要关注海量数据集中最重要的信息<sup>[1]</sup>。数据挖掘中分类这块获得了很大的关注度, 基于分析一系列样本数据形成分类的模型来预测未来数据对象的所属类别。

### 1 前言

一般来说, 在分类中, 规则以 IF...THEN 的形式

来表示知识: IF (条件) THEN (类别)。在这里, 规则提取的主要目的是找出隐藏的知识, 并且以可理解的方式来解释它。分类方法的代表方法便是神经网络, 神经网络在预测、从经验中学习、从先前样本中泛化的能力是十分优秀的, 这使得它们非常适合分类问题的应用。它能够获得很高的分类准确性, 然而它的一个重要的缺陷是其解释能力非常有限<sup>[2]</sup>。这主要因为神经网络中的知识是分散在其激活函数和神经元的联接上<sup>[3]</sup>。为了使人们能够容易的理解神经网络中隐藏的知识, 近年来研究人员试图发展智能的算法, 以便

① 收稿时间:2010-10-31;收到修改稿时间:2010-12-12

从训练好的神经网络中提取知识，如遗传算法 (Elalfi 等<sup>[4]</sup>)、进化算法、蚁群算法、支持向量机、粗糙集、人工免疫系统 (Humar Kahramanli 等<sup>[5]</sup>)。

由于蚁群算法与神经网络相结合的研究相对较少，而 Lale Ozbakir<sup>[6]</sup>等人在蚁群算法与神经网络相结合进行分类工作上取得了很好的进展。因此本文将研究的重点放在神经网络与改良的蚁群算法结合，取神经网络强大的自我训练能力和蚁群算法解决优化问题的能力来达到更好的理解神经网络，获取其中的知识，产生分类规则。

## 2 相关知识介绍

### 2.1 神经网络简介

神经网络是一个从生物神经系统启发而来的信息处理机制，由一系列的内部相互联系的神经元的工作来解决特定的问题。内部神经元之间的联接决定了神经网络的结构。学习算法主导修正联接或者训练来实现所需的神经网络的行为能力<sup>[7]</sup>。由于神经网络从复杂的，不精确的数据中追溯出意义的优秀能力，神经网络能够用于提取出那些太过复杂以至于不被人所注意的模式。但由于神经网络的黑箱特性，解释神经网络的学习过程是很困难的。因此近年来研究人员关注于开发神经网络使人类能够理解的表达方式。本文便是通过改良的蚁群算法对训练后的神经网络进行分类规则提取。

### 2.2 蚁群算法简介

ACO(Ant Colony Optimization)蚁群算法是 Dorigo 为解决组合优化问题而提出的。ACO 算法是由蚂蚁觅食的行为而启发的。在大多数的蚂蚁群体中，最初蚂蚁以一定的行为搜索它们巢穴附近的区域。只要有蚂蚁找到了食物，则在评估完食物的数量和质量以后，蚂蚁便带一部分食物会巢穴。在回去的路上，蚂蚁根据食物的数量和质量这些信息释放一些信息素，以便其它蚂蚁也能够找到该食物。路上的释放信息素以一定的概率帮助蚂蚁找到食物。蚂蚁通过信息素来进行这种非直接的交流，使它们找到食物和巢穴的最短路径<sup>[8]</sup>。

本文将采用一个不同的蚁群算法，从 TACO (Touring Ant Colony Optimization)中改良而来。TACO 算法由 Hiroyasu 等<sup>[9]</sup>首先提出，为了处理连续型变量优化问题。与普通的蚁群算法不同的是，TACO 算法

中每个方案以二进制串的形式表示。蚂蚁试着选择每个二进制位的值是 0 或 1。算法的概念如图 1 所示，算法的具体内容和其改进将在后面提出。

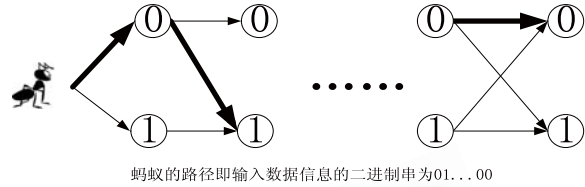


图 1 TACO 算法二进制串化

## 3 基于改良蚁群算法的神经网络分类规则提取算法

算法通过改良蚁群算法来对训练后的神经网络进行分类规则的提取。主要步骤如下：

- ① 首先要对分类的数据进行预处理，也就是以二进制串的形式来表示离散型数据或者离散化后的连续性数据，以便作为神经网络的输入数据。
- ② 进行神经网络的训练，数据预处理后的二进制串为输入，分类类别为输出。
- ③ 用改良后的 TACO 和训练好的神经网络相结合进行规则提取。
- ④ 对规则进行评估，选取满足要求的规则，精简或者合并规则形成一套完整的规则集。

由于 TACO 对每个神经网络输入属性进行二进制串化，导致属性过多的情况下，方案将变的很长，影响了性能，并且由于同一属性只能取一个值，即每个属性的二进制串只能有一位为 1，导致算法中蚂蚁所走的路径需要保存大量信息，即需要记录同一个属性中是否有一位为 1，从而其他值只能为 0，算法复杂而不直观。以下是我对该算法的改进，如图 2。

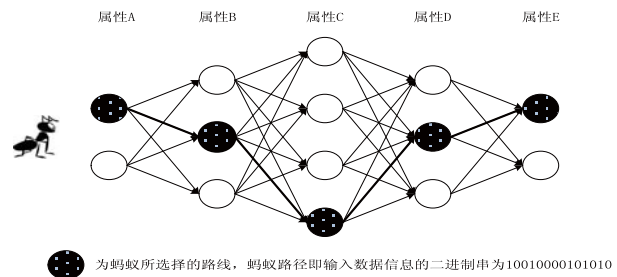


图 2 改良 TACO 算法二进制串化

由于 TACO 分类规则的提取是对神经网络规则提

取进行的简单化、自适应化和可理解性的研究，除了预处理后的数据二进制串既作为神经网络输入，又是 TACO 的方案，其与神经网络的数据训练能力和准确率等是可以脱离开来的，并且对神经网络训练的研究比较成熟，并获得了很好的分类准确率，因此本文的重点放在规则提取这块，以便其能够和已经存在的优秀的神经网络的分类算法相结合。

### 3.1 数据预处理

表 1 是一个潜在顾客分类的例子，数据包含离散属性，如性别；连续属性，如年龄，月收入。首先对属性进行分类，对连续型和离散型属性进行划分，如月收入可划分为一般收入，中等收入，高收入，然后进行二进制编码，这样每一条顾客信息均被编码为一个二进制字符串。如表 1 所示，张三的二进制码为 10100100。而提取出的规则形式可能如下：IF 性别是男性 AND 年龄是青年 AND 月收入是中等收入，THEN 潜在顾客。在数据预处理中对连续型属性的划分，直接影响到分类及规则提取。本文使用 UCI 的数据集，并使用数据挖掘软件 weka 对数据集进行分析，从而能够更加直观地进行属性划分。

表 1 潜在顾客分类的例子

属性	性别		年龄			月收入		
	男	女	18-30	31-50	51-	小于 3K	3K-10K	大于 10k
二进制	1	0	100	010	001	100	010	001
张三	男		25			2800		

### 3.2 改良 TACO 算法产生规则

改良 TACO 规则提取算法的参数如下：

N: 蚂蚁数量；T: 迭代次数； $\rho$ : 信息素挥发率；f: 频率因子；Q: 正变量； $P_{ij}$ : i-j 的子路径概率； $f_{ij}$ : i-j 的子路径频率； $\tau_{ij}$ : 子路径 i-j 上的信息素数量； $\Delta\tau_{ij}^k$ : 在时间 t 到 t+1 之间由第 k 只蚂蚁附加到子路径 i-j 上的信息素数量； $\Delta\tau_{ij}(t, t+1)$ : 在时间 t 到 t+1 之间附加到子路径 i-j 上的信息素总量； $\tau_{ij}(t+1)$ : 在时间点 t+1 时，子路径 i-j 上的信息素数量。注：i-j 指的如图 2 中两个属性间所选的路线。

质量函数：用于评价规则的质量。 $S_e$  指正例被正确分类， $S_p$  指反例被正确分类。tp 和 tn 指正例和反例

被正确分类，fp 和 fn 指正例和反例被错误分类。

$$Quality = S_e * S_p = \frac{tp}{tp + fn} * \frac{tn}{tn + fp} \quad (1)$$

产生规则：

每个质量函数为正值的路径被存储为候选规则用于规则精简过程。当 Quality 最大化的时候，便可知道  $x_i=0$  或 1，即二进制字符串的值，作为一条规则。编码和初始化蚁群：二进制编码在上面已经介绍了，初始化蚁群为随机产生。

产生蚂蚁的路径：

蚁群算法容易导致局部收敛，为了防止局部收敛现象产生，需要对每条路径上给蚂蚁一个选择其他并非最优路径的概率。如公式 2 所示

$$p_{ij}(t) = \begin{cases} 1 \\ \frac{\tau_{ij}}{\tau_{i1} + \dots + \tau_{in}} \end{cases} \quad (2)$$

假如  $f * f_{ij} < (f_{i1} + \dots + f_{in})$ ，则选择路径 ij 不然选

择路径 ij 的概率为  $\frac{\tau_{ij}}{\tau_{i1} + \dots + \tau_{in}}$

更新信息素和频率值：

信息素和频率由公式 3, 4 同时计算。 $F_k$  是蚂蚁 k 完成路径后计算的目标函数值，即上面的 Quality 质量函数，Q 是一个固定常量。

$$\Delta\tau_{ij}^k(t, t+1) = \begin{cases} Q * F_k \\ 0 \end{cases} \quad (3)$$

当有蚂蚁通过子路径 ij，则信息素增加为  $Q * F_k$

$$\Delta F_{ij}^k(t, t+1) = \begin{cases} 1 \\ 0 \end{cases} \quad (4)$$

当有蚂蚁通过子路径 i-j，则频率加 1

当所有蚂蚁完成了搜索过程并产生了路径，则在时间 t 到 t+1 之间附加到子路径信息素和频率值的计算公式为 5, 6:

$$\Delta\tau_{ij}(t, t+1) = \sum_{k=1}^N (t, t+1) \quad (5)$$

$$\Delta F_{ij}(t, t+1) = \sum_{k=1}^N (t, t+1) \quad (6)$$

同一条子路径在时刻 t+1 时的信息素和频率值由公式 7, 8 更新。

$$\tau_{ij}(t+1) = \Delta\tau_{ij}(t, t+1) + \rho\tau_{ij}(t) \quad (7)$$

$$F_{ij}(t+1) = \Delta F_{ij}(t, t+1) + F_{ij}(t) \quad (8)$$

算法结束条件以及规则精简：

算法执行到一个特定的迭代值，即参数 T。算法结束后，将得到一个候选方案的列表，这些方案有可能成为规则。根据候选方案的质量函数值来进行规则精简，函数值最大的候选方案即为最合适的方案，加入到一个规则列表中，从训练数据中移除满足该方案的数据；将该方案作为规则保存在规则列表中；直到所有的训练数据都在规则列表中有对应的规则或者候选方案全部处理完。

#### 4 实验结果及分析

实验采用鸢尾花(Iris)数据集,从加州大学厄文分校(UCI)的机器学习库中得到。鸢尾花数据集包含 150 条信息,三个鸢尾花种:Setosa、Versicolour 和 Virginica 各 50 条,包含 4 个连续属性:萼片长度,萼片宽度,花瓣长度,花瓣宽度。其取值为[4.3,6.90], [2.0,4.4], [1.0,6.9],[0.1,2.5]。

数据预处理用 matlab 实现,首先将 weka 数据挖掘软件的 arff 格式转化为 xls 格式,并读入 matlab,根据事先对 Iris 的 4 个属性进行分类的条件,进行数据的二进制串化。

而对属性的分类极大地影响到规则的提取,需要尽可能的将相似的数据归为一类,如果分类范围过大,或者所分类别较少,则提取的规则可能不具有典型性,比如将萼片长度分为长和短两类不如分为长,中,短三类能够提出更好的规则。而分类范围以簇类的形式进行分类,尽可能的将属性值较靠近的分为一类,能够提取更好的规则。

通过 weka 软件的对数据集直观的数据分析加上对 Iris 数据集 150 个数据项的分析,对 Iris 数据集属性如下表 2。

表 2 Iris 属性分类

萼片 长度	x≤5.45		5.33<x≤7		x>7	
	100	C1	010	C2	001	C3
萼片 宽度	x≤2.9		2.9<x≤3.4		If(x>3.4)	
	100	C1	010	C2	001	C3
花瓣 长度	x≤2.45		2.45<x≤4.75		x>4.75	
	100	C1	010	C2	001	C3
花瓣 宽度	x≤0.8		0.8<x≤1.75		x>1.75	
	100	C1	010	C2	001	C3

实验的数据预处理采用 Matlab 实现,规则提取部分则由 C++语言实现。实验结果如下,提取出规则三条:

If 萼片长度为 C1,萼片宽度为 C2 或 C3,花瓣长度为 C1,花瓣宽度为 C1,Then 类型为 Setosa。Setosa50 条数据中有 7 条不符合该规则,规则符合度 86%。

If 萼片长度为 C2,萼片宽度为 C1 或 C2,花瓣长度为 C2,花瓣宽度为 C2,Then 类型为 Versicolour。Versicolour50 条数据中有 11 条数据不符合该规则,规则符合度为 78%。

If 萼片长度为 C2 或 C3,萼片宽度为 C1 或 C2,花瓣长度为 C3,花瓣宽度为 C3,Then 类型为 Virginica。Virginica 中有 5 条数据不符合该规则,规则符合度为 90%。

实现结果表明:提取的规则能够很大程度上反映了鸢尾花的所属类别之间在各个属性上的差异,一般正常情况下的花种都能与其所属花种相对应,也即改良蚁群算法很大程度上是寻找不同种类上的典型差异。改良蚁群算法其算法特点便是通过数据预处理使各个属性之间具有更细的差异,使每个对分类有影响的属性在更小的范围内影响到分类的结果。而且改良蚁群算法能够从数据集中提取较少规则,这样所提取的规则所覆盖的范围更精确,而不是找出很空泛而无用的规则。与其它基于规则的分类方法相比,如 DOEA<sup>[10]</sup>平均规则数 4 条,GARC<sup>[11]</sup>平均规则数 7 条,但其前提是能够对数据属性进行合理的分类以及筛选,比如当属性量过大的时候,对每个属性分类,进行范围区间选择关系到蚁群算法的规则提取性能,所以对属性的合理分类和筛选是一个需要考虑的问题;并且由于属性分类过多会导致神经网络输入节点过多,影响到训练速度,所以可以在对数据属性的预处理中加入计算属性对神经网络的影响因子的大小进行属性筛选,从而尽可能地降低其影响。

本文提出了一种新颖的改良蚁群算法用于神经网络分类规则的提取,通过对 Iris 数据集的分类规则提取应用,可以看出改良蚁群算法优秀的规则提取能力,能够获得更精确规则,将其与优秀神经网络分类算法相结合,能够得到很好的实用性。

## 参考文献

- 1 Hand DJ, Mannila H, Smyth P. Principles of Data Mining (Adaptive Computation and Machine Learning), MIT Press, 2001.
  - 2 Andrews R, Diederich J, Tickle AB, A survey, critique of techniques for extracting rules from trained artificial neural networks, Knowledge Based Systems, 1995,8(6): 373–389.
  - 3 Bologna G. Is the worth generating rules from neural network ensembles. Journal of Applied Logic, 2004,(2): 325–348.
  - 4 Elalfi AE, Haqueeb R, Elalami ME. Extracting rules from trained neural network using GA for managing E-business. Applied Soft Computing, 2004,(4):65–77.
  - 5 Kahramanli H, Allahverdi N. Rule extraction from trained adaptive neural networks using artificial immune systems. Expert Systems with Applications, 2009,(36): 1513–1522.
  - 6 Lale Özbakir a, Adil Baykasoglu b, Sinem Kulluk a, Hüseyin Yapici c. An ant colony based algorithm for rule extraction from trained neural networks. Expert Systems with Applications, 2009. doi:10.1016/j.eswa.2009.04.058
  - 7 Diego A. Computational Intelligence: For Engineering and Manufacturing. Boston, Springer US, MA: USA, 2007.
  - 8 Dorigo M, Maniezzo V, Colomi A. Positive feedback as a search strategy. Technical Report N.91-016 Politecnico di Milano, 1991.
  - 9 Hiroyasu T, Miki M, Ono Y, Minami, Y. Ant colony for continuous functions. The Science and Engineering, Doshisha University, 2000.
  - 10 Tan C, Yu Q, Ang JH. A dual-objective evolutionary algorithm for rules extraction in data mining. Computational Optimization and Applications, 2006, 34: 273–294.
  - 11 Chen T, Hsu TA. Gas based approach for mining breast cancer pattern. Expert Systems with Applications, 2006,30: 674–681.
- 
- (上接第 110 页)
- 4 Azzedin F, Maheswaran M. Evolving and Managing Trust in Grid Computing Systems. IEEE Canadian Conference on Electrical & Computer Engineering (CCECE'02). May 2002. 1424–1429.
  - 5 Alunkal B, Veljkovic L, Laszewski GV, et al. Reputation-based Grid resource selection. Proc. of the Workshop on Adaptive GridMiddleware. New Orleans, 2003.
  - 6 王莉苹,杨寿保.网格环境中的一种信任模型.计算机工程与应用,2004,40(23):50–53.
  - 7 王珊,高迎,程涛远,等.服务网格环境下基于行为的双层信任模型的研究.计算机应用,2005,9:1974–1991.
  - 8 王东安,徐浩,南凯,等.基于推荐的网格计算的信任模型计算机应用研究.计算机应用研究,2006,2:96–98.
  - 9 李文娟,王晓东,傅仰歌,傅志祥.几种网格信任模型的研究.福州大学学报(自然科学版),2006,34(2):190–193.