

基于决策树算法的遥感图像分类研究^{*}

罗来平¹, 宫辉力², 刘先林²

(1. 北京城市学院 人工智能研究所, 北京 100083; 2. 首都师范大学 资源环境与旅游学院 资源环境与地理信息系统 北京市重点实验室, 北京 100037)

摘 要: 针对传统分类方法在处理空间特征分布极为复杂的数据时效果不佳的缺点, 结合分层思想的树分类技术, 对广泛用于数据挖掘模型中的 CART 决策树算法进行改进, 提出了一种基于人机交互的决策树算法, 将其应用到遥感图像自动分类中, 具有很好的弹性和鲁棒性, 且分类结构简单明了, 达到了更好的分类效果。以 VC++ 6.0 作为开发工具, 定义了一种特殊的数据结构, 实现了该分类系统。实践表明, 该系统具有很好的稳定性和交互性, 实用性较强。

关键词: 决策树; 算法; 图像分类; 遥感; VC++

中图分类号: TP391 **文献标识码:** A **文章编号:** 1001-3695(2007)01-0207-03

Research and Implementation of Classification in Remote Sensing Image Based on Decision Tree Algorithm

LUO Laiping¹, GONG Hui-li², LIU Xian-lin²

(1. Artificial Intelligence Institute, Beijing City University, Beijing 100083, China; 2. Key Laboratory of Resource Environment & GIS, College of Resource Environment & Tourism, Capital Normal University, Beijing 100037, China)

Abstract: In allusion to the traditional classification's shortcoming of low effect when dealing with remote sensing data that has complex spatial-character distributing, combining with tree-classification technology which using delaminating method, Classification and Regression Tree (CART) algorithm is improved, which is used largely in data mining. And an algorithm based on human-computer interaction is put forward, which is applied to remote sensing image classification with better flexibility and robust, also with simple and clear class structure, having achieved better classified effect. And a special data structure is designed to realize this classifying system in VC++. Practice show that this system has good stability, alternateness and strong practicability.

Key words: Decision Tree; Algorithm; Image Classification; Remote Sensing; VC++

在遥感技术的研究中,通过遥感图像判读识别各种目标是遥感技术发展的一个重要环节^[1],无论是专业信息获取、动态变化预测,还是专题地图制作和遥感数据库的建立都离不开分类。早期的分类技术是目视解译,其时效性、可重复性差,解译结果也因人而异,很难进行比较和转换。近年来随着计算机技术的飞速发展,计算机识别自动分类已逐渐代替了早期的分类技术,成为遥感应用的一个重要组成部分,也是当前遥感发展的前沿^[2]。目前,遥感图像自动分类主要利用统计模式方法^[1],可分为非监督和监督两类。由于遥感图像往往存在异物同谱和同物异谱的现象,用传统的统计模式方法分类的效果不甚理想,因而人们不断研究新的分类方法。例如模拟人脑思维方式提出的人工神经网络分类^[3];针对地物特征和模糊性提出的模糊聚类分类^[4];模式识别与人工智能技术相结合的专家系统分类^[1];计算混合像元内各典型地物所占面积比例的混合像元分解法^[5];按照一定的知识规则,采用分层思想的树分类^[6]等技术。

虽然这些方法在一定程度上提高了分类精度,达到了更好的效果,但是它们仍是以遥感图像的光谱特征为基础的。在实践中,由于各种因素的影响,如地面起伏对地物光谱反射强度的影响^[7];像素的分辨率低而造成的混类像素影响;分类中只考虑单个像素光谱特征而未考虑相邻像素类属的相关性及结构特征等因素,都会使常规的计算机分类效果不够理想。因此需要采取一些辅助的处理措施^[1],如引入地面高程、坡度、坡向信息等,以设法改善分类效果。但目前的遥感图像计算机自动分类技术均不能很好地引入这些辅助信息^[1]。本文在采用分层思想的树分类的基础上,对广泛用于数据挖掘模型中的决策树算法进行改进,并以 VC++ 6.0 作为开发工具,将其应用到图像自动分类中,能够很好地解决这个难题,达到了更好的分类效果。

1 决策树算法的图像分类研究及实现

1.1 常用的决策树算法简介

决策树^[8,9]是一个类似流程图的树型结构,其中树的每个内部节点代表对一个属性的测试,其分支代表测试的每个结

收稿日期: 2005-10-18; 修返日期: 2006-11-20
基金项目: 国家“863”计划资助项目(2003AA135010)

果,而树的每个叶子节点代表一个类别,树的最高层节点就是根节点,是整个决策树的开始。其算法的早期版本可以追溯到20世纪60年代^[10]。后来特别是最近几年,它被广泛应用到许多需要分类识别的领域,如科学实验、医疗诊断、信贷审核、商业预测、案件侦破等。这类算法无须相关领域知识,且相对于基于模糊理论的分类方法,具有更高的分类准确率和更快的处理速度^[11]。

在很多领域特别是数据挖掘中,决策树是一种经常要用到的技术^[12],它可以用于分析数据,也可以用来作预测,常用的算法有D3, CART, C4.5等。

(1) D3算法是最有影响和最早的决策树算法之一^[10],其建立在推理系统和概念学习系统的基础上,但它是非递增学习算法。每当一个或数个新例子进来,就必须重新执行一次该算法,把新来的例子和以前旧的全部例子集合变成决策树,因此效率非常低。而且它是基于单变量的,难以表达复杂概念,抗噪性差。

(2) C4.5是D3的改进版本^[13]。它主要在以下几个方面对D3作了改进:缺省值的预测属性仍可用,提出了修剪思想,可以进行规则推导。

(3) CART (Classification and Regression Tree, 分类回归树)^[14]是一种数据勘测和预测算法。它用一种非常简单的方法来选择问题,即将每个问题均试一次,然后挑出最好的一个,用它把数据分成更有序的两个分割,再对新的分割分别提出所有可能的问题。因此该算法得到的决策树每个节点有两个分支,即二叉树。

1.2 改进决策树生成算法

为了更好地将决策树算法应用于遥感图像分类中,本文在CART算法的基础上作了以下改进:

(1) 常用的决策树算法均由用户提供训练样本集,计算机执行算法生成决策树,整个过程由计算机自动完成,不需要任何人工干预。由于遥感图像训练集的属性取值太多,而有些取值是用不到的,需要把这些用不到的取值过滤掉,否则会影响整颗树的质量。因此本文在决策树的生成过程中引入人机交互技术,将用户的先验知识用于决策树的生成过程中,使得生成的决策树更加合理可信。基于人机交互的决策树方法由用户与计算机相互交互,共同完成,故要求计算机提供可视化环境和工具(图1),友好的界面方便用户输入先验知识。

(2) 建立一棵树首先需要选择一个属性作为根节点,然后将该属性的每一个可能值作为一个分支;再在每个分支所剩余的属性中找出一个属性作为该分支的下一个节点,如此循环到所有属性均被选用为止。常用的决策树算法均采用信息熵作为选择标准。由于遥感图像中的噪音比较多,而基于信息熵属性选择标准往往抗干扰能力不强且以对数计算为累加计算的计算量较大,故本文采用了一种新型的属性选择标准——属性重要性^[15]来提高属性选择的效率。该方法是用训练值的变化而引起输出变化的累加值作为衡量属性重要性的标准,即对于某个属性,如果训练值的变化而引起的输出变化越大,说明该属性就越重要。可用式(1)表示为

$$C(K) = |x(i, k) - x(j, k)| \times \text{sign}|y(i) - y(j)| (i, j) \quad (1)$$

式中: $C(K)$ 表示第 k 个属性的输入/输出关联值; $x(i, k)$, $x(j, k)$ 表示第 i , j 个样本的第 k 个条件属性值; $y(i)$, $y(j)$ 表示第 i , j 个样本的决策属性值; $\text{sign}(x)$ 表示符号函数。

1.3 一种自定义的数据结构

数据结构的设计在程序设计中很关键,一个好的数据结构可以让算法更加精练,大大提高开发效率。本文采用的算法是在CART算法的基础上改进过来的,因此生成的也是一棵二叉树。但由于在生成算法中要频繁地查找和调整决策树(添加或删除子树),传统的二叉树结构在速度上不能满足算法的需求,故笔者在设计系统时创新了一种类二叉树的结构(图2)。在树的每一层都设置了一个头节点,且树的每个节点只有指向父节点和左右兄弟节点的指针。

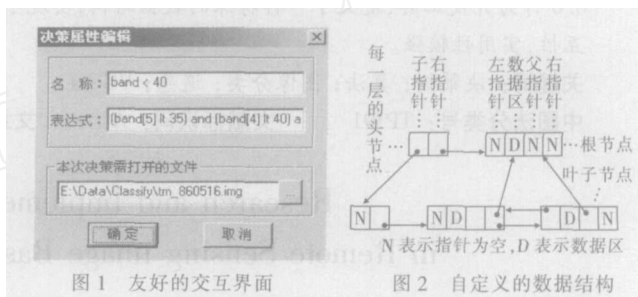


图1 友好的交互界面

图2 自定义的数据结构

```
class TreeNode // 树节点的结构
{
public:
    TreeNode Info m_NodeInfo; // 这个节点的数据区
    TreeNode * m_ParentNode; // 父节点指针
    TreeNode * m_LBtotherNode; // 左兄弟节点指针
    TreeNode * m_RBtotherNode; // 右兄弟节点指针
    TreeNode(position newpos, CString newQorA, CString newExpression);
    TreeNode();
    virtual ~TreeNode();
};

struct level_node // 每一层的头节点结构
{
    TreeNode * level_treenode; // 指向这一层的第一个树节点
    level_node * next_level; // 指向下一层的头节点
};
```

1.4 系统实现

图3为系统架构图。其中 D 是训练集合, A 为分类属性集合。另有测试数据集合 T 用来评估生成决策树的误差。整个过程分为两个部分进行: 决策树构造, 也称学习过程, 主要工作是输入训练集, 采用改进的决策树生成算法生成决策树, 并作好分类前的预备工作, 即提取分类规则; 决策树预测, 也称分类过程, 主要工作是应用分类规则进行分类, 并根据测试集, 计算出分类误差, 误差较大的对决策树作裁剪算法。

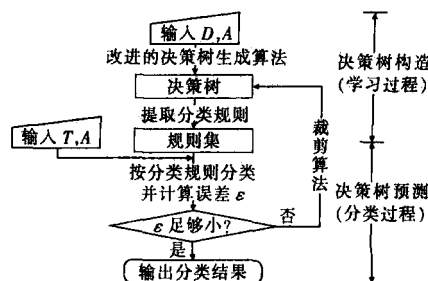


图3 系统架构图

误差较小的就输出其分类结果,分类结果有两种:

(1)将分类结果写入新的遥感图像文件(如 MG 文件)中;

(2)可视化的树结构,在树的叶子节点保存着每一类的属性。

该系统以 VC++6.0 为开发工具,采用面向对象的思想设计与开发。

2 应用实例

本文以某农村 TM1-7 波段(图 4)作为数据源,利用本系统进行决策树分类。图 5 为分类后的图像,图 6 为分类后决策树分类结果图。为了对分类精度进行评价,本文将原始图像与分类后的图像进行对比,选取了大量的样本数据,建立了该分类方法的混淆矩阵,并计算出总体分类精度为 86.9%, Kappa 系数可达 0.85。通过混淆矩阵发现,除了建筑用地外,对于其他地物可以达到 90.3% 以上的分类精度。经检验发现,其误判像元主要是位于耕地中。这主要是因为建筑用地与耕地形状都非常相似,且紧密分布在一起,从而在分类中被误判。

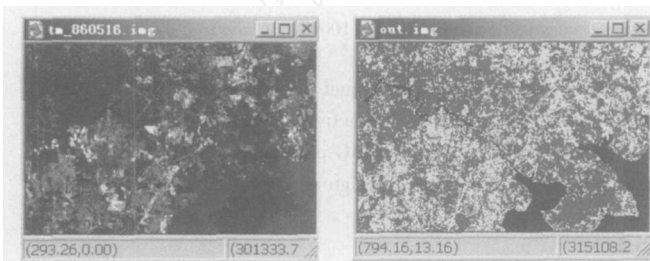


图 4 原始图像灰度图

图 5 决策树分类后的图像

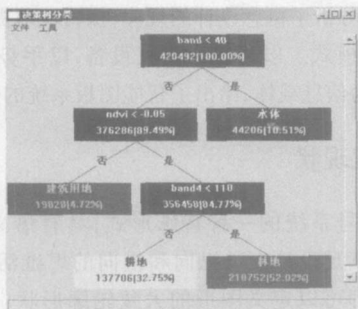


图 6 决策树分类结果图——二叉树

3 总结与展望

在遥感技术飞速发展的今天,人们更多地应用多源遥感数据、多波段图像数据、高光谱数据,如何充分利用这些数据,又不至于使计算速度降低,传统分类方法就显得力不从心了,尤其是面对空间特征分布极为复杂的数据时,更显示出它们的不足。通过本文研究发现:决策树算法对于输入数据的空间特征和分类标志具有更好的弹性和鲁棒性,它用于遥感数据分类的优势主要在于对数字图像数据特征空间的分割上,其分类结构简单明了,尤其是二叉树结构的单一决策树结构十分容易解释。因此,当遥感图像数据特征的空间分布很复杂,或者源数据各维具有不同的统计分布和尺度时,基于决策树算法的分类方法能够获得较为理想的分类结果。然而,决策树算法应用于

遥感图像分类尚处于探索阶段,大部分都是理论上针对某一具体应用的探讨,目前研究成果比较少^[16],还没有较为通用的软件诞生。本文则在探索理论的基础上对开发该分类系统作出了尝试,并且该分类系统已投入遥感应用行业中,实践发现该系统交互性和稳定性很好,能适应一般的需求。进一步的工作就是将决策树算法与其他技术,如神经网络相结合应用到该系统中,以期获取更高的分类精度和更好的分类效率。

参考文献:

- [1] 汤国安,等. 遥感数字图像处理 [M]. 北京:科学出版社,2004.
- [2] 张仁华. 实验遥感模型及地面基础 [M]. 北京:科学出版社,1996.
- [3] Ito Y, Omatu S. Extended LVQ Neural Network Approach to Land Cover Mapping [J]. IEEE Transactions on Geoscience and Remote Sensing, 1999, 37 (1): 313-316.
- [4] Paola J D, Schowengerdt R A. A Detailed Comparison of Back-propagation Neural Network and Maximum-likelihood Classifiers for Urban Land Use Classification [J]. IEEE Transactions on Geosciences Remote Sensing, 1995, 33 (4): 981-996.
- [5] Fridel M A, Brodley C E, Strahler A H. Maximizing Land Cover Classification Accuracies Produced by Decision Trees at Continental to Global Scales [J]. IEEE Transactions on Geosciences and Remote Sensing, 1999, 37 (2): 969-979.
- [6] Friedl M A, Brodeley C E. Decision Tree Classification of Land Cover from Remotely Sensing Data [J]. Remote Sensing of Environment, 1997, (61): 399-409.
- [7] 李彤,等. 采用决策树分类技术对北京市土地覆盖现状进行研究 [J]. 遥感技术与应用, 2004, 19 (6): 485-487.
- [8] Jiawei Han, Micheline Kamber. Data Mining: Concepts and Techniques [M]. 北京:高等教育出版社,2002.
- [9] Alex Berson, Stephen Smith, Kurt Thearling. Building Data Mining Application for CRM [M]. 北京:人民邮电出版社,2001.
- [10] 戴南. 基于决策树的分类方法研究 [D]. 南京师范大学硕士学位论文, 2003.
- [11] 李宁,等. 决策树算法及其常见问题的解决 [J]. 计算机与数字工程, 2005, 33 (3): 60-64.
- [12] 孙雪莲. 数据挖掘中分类算法研究 [D]. 吉林大学硕士学位论文, 2002.
- [13] 曹叶虹. 结合粗糙集理论的决策树技术的研究 [D]. 华南理工大学硕士学位论文, 2002.
- [14] 秦欢,等. 基于决策树 CART 的中文文语转换系统语音合成单元的预选 [J]. 微型电脑应用, 2004, 20 (5): 5-6.
- [15] 倪春鹏,等. 一种新型决策树属性选择标准 [J]. 武汉科技大学学报(自然科学版), 2004, 27 (4): 437-439.
- [16] 刘小平,等. 像元信息分解和决策树相结合的图像分类方法 [J]. 地理与地理信息科学, 2004, 20 (6): 35.

作者简介:

罗来平 (1982-), 男, 江西人, 助教, 工学硕士, 主要研究方向为遥感与三维地理信息系统; 宫辉力 (1956-), 男, 长春人, 教授, 博导, 博士, 主要研究方向为地理信息系统和遥感技术应用; 刘先林 (1939-), 男, 广西人, 院士, 博导, 主要研究方向为三维获取与应用。