

个人资料



moodytong

访问： 279692次

积分： 3930

等级： BLOG > 5

排名： 第4349名

原创： 115篇 转载： 24篇

译文： 6篇 评论： 50条

文章搜索

文章分类

Java/python (34)

c/c++ (12)

linux (4)

mplayer (2)

云计算与SOA (14)

图 (6)

学术科技 (20)

数据结构与算法 (15)

日常问题 (7)

琐事杂谈 (5)

复杂网络&网络科学 (12)

统计分析与数据挖掘 (17)

阅读排行

python encode和decode (22703)

Python学习之字典详解 (17786)

利用matlab进行简单的贝 (16133)

networkx使用笔记(一)之: (8607)

networkx使用笔记(二)之: (7930)

networkx使用笔记(三)之: (7319)

机器学习经典算法8-树回归 (6535)

机器学习经典算法7-线性 (6447)

学院APP首次下载，可得50C币！ 欢迎来帮助开源“进步” 当讲师？爱学习？投票攒课吧 CSDN 2015博客之星评选结果公布

机器学习经典算法8-树回归

2013-08-19 09:18

6536人阅读

评论(6)

收藏

举报

分类： 统计分析与数据挖掘 (16)

版权声明：本文为博主原创文章，未经博主允许不得转载。

目录(?)

[+]

1. 简单介绍

线性回归方法可以有有效的拟合所有样本点（局部加权线性回归除外）。当数据拥有众多特征并且特征之间关系十分复杂时，构建全局模型的想法一个是困难一个是笨拙。此外，实际中很多问题为非线性的，例如常见的分段函数，不可能用全局线性模型来进行拟合。

树回归将数据集切分成多份易建模的数据，然后利用线性回归进行建模和拟合。这里介绍较为经典的树回归

CART(classification and regression trees, 分类回归树)算法。

2. 分类回归树基本流程

构建树：

1. 找到[最佳待切分特征]
2. 若不能再切分，则将该节点存为[叶子节点]并返回
3. 按照最佳待切分特征将数据集切分成左右子树（这里为了方便，假设大于特征值则为左，小于则归为右）
4. 对左子树进行[构建树]
5. 对右子树进行[构建树]

最佳待切分特征：

1. 遍历特征
 - 1.1 遍历特征所有特征值
 - 1.1.1 计算按该特征值进行数据集切分的[误差]
2. 选择误差最小的特征及其相应值作为最佳待切分特征并返回

基于回归树的预测：

1. 判断当前回归树是否为叶子节点，如果是则[预测]，如果不是则执行2
2. 将测试数据相应特征上的特征值与当前回归树进行比较，如果测试数据特征值大，则判别当前回归树的左子树是否为叶子节点，如果不是叶子节点则进行[基于回归树的预测]，如果是叶子节点，则[预测]；反之，判别当前回归树的右子树是否为叶子节点，如果不是叶子节点则进行[基于回归树的预测]，如果是叶子节点，则[预测]

3. 分类回归树的实践说明

误差、叶子节点和预测三者有相关的关联关系，一种相对简单的是误差采用的是y值均方差，叶子节点相应的建立为该节点下所有样本的y值平均值，预测的时候根据判断返回该叶子节点下y值平均值即可。

在进行最佳待切分特征选取的时候，一般还有两个参数，一个是允许的误差下降值，一个是切分最小样本数。对于允许误差下降值，在实际过程中，需要在分割之后其误差减少应该至少大于该bound；对于切分最小样本数，也就是说切分后的子树中包含的样本数应该多于该bound。其实这两种策略都是为了避免过拟合。

4 树剪枝

密...

关于机器学习课程的小记

longyou1243: 资料的连接不能用了, 楼主更新一下连接可否?

机器学习经典算法10-Apriori

拾毅者: 36和37行还是不太明白, 看《机器学习实战》这两行有点不明白, 能解释下么
L1=list(Lk) L...

社交标签技术的研究

eBruce: 博主对“标签”理解太深了, 彻底佩服。“标签”几乎可以在所有平台看到, 只是形态上有些区别。最近“图片+...

机器学习经典算法7-线性回归

月光下的夜曲: 代码有错误, 害人啊

networkx使用笔记(三)之好汉篇V

moodytong: @novalist在前面要申明下import numpy as np注意安装numpy库。此外, 有...

```
51. for featIndex in range(n-1):
52.     for splitVal in set(dataSet[:,featIndex]):
53.         mat0, mat1 = binSplitDataSet(dataSet, featIndex, splitVal)
54.         if (shape(mat0)[0]<tolN) or (shape(mat1)[0]<tolN):
55.             continue
56.         newS = errType(mat0)+errType(mat1)
57.         if newS < bestS:
58.             bestIndex = featIndex
59.             bestValue = splitVal
60.             bestS = newS
61. if (S-bestS)<tolS:
62.     return None, leafType(dataSet)
63. mat0, mat1 = binSplitDataSet(dataSet, bestIndex, bestValue)
64. if (shape(mat0)[0]<tolN) or (shape(mat1)[0]<tolN):
65.     print "Not enough nums"
66.     return None, leafType(dataSet)
67. return bestIndex, bestValue
68. def binSplitDataSet(dataSet, feature, value):
69.     mat0 = dataSet[nonzero(dataSet[:, feature]>value)[0],:] [0]
70.     mat1 = dataSet[nonzero(dataSet[:, feature]<=value)[0],:] [0]
71.     return mat0, mat1
72. def createTree(dataSet, leafType=regLeaf, errType=regErr, ops=(1,4)):
73.     feat, val = chooseBestSplit(dataSet, leafType, errType, ops)
74.     if feat == None:
75.         return val
76.     retTree={}
77.     retTree['spInd'] = feat
78.     retTree['spVal'] = val
79.     lSet, rSet = binSplitDataSet(dataSet, feat, val)
80.     retTree['left']=createTree(lSet, leafType, errType, ops)
81.     retTree['right']=createTree(rSet, leafType, errType, ops)
82.     return retTree
83. def isTree(obj):
84.     return (type(obj).__name__=='dict')
85. def getMean(tree):
86.     if isTree(tree['right']):
87.         tree['right'] = getMean(tree['right'])
88.     if isTree(tree['left']):
89.         tree['left'] = getMean(tree['left'])
90.     return (tree['left']+tree['right'])/2.0
91. def prune(tree, testData):
92.     if shape(testData)[0] == 0:
93.         return getMean(tree)
94.     if (isTree(tree['right']) or isTree(tree['left'])):
95.         lSet, rSet = binSplitDataSet(testData, tree['spInd'], tree['spVal'])
96.         if isTree(tree['left']):
97.             tree['left']=prune(tree['left'], lSet)
98.         if isTree(tree['right']):
99.             tree['right']=prune(tree['right'], rSet)
100.         if not isTree(tree['right']) and not isTree(tree['left']):
101.             lSet, rSet = binSplitDataSet(testData, tree['spInd'], tree['spVal'])
102.             errorNoMerge = sum(power(lSet[:,-1]-tree['left'],2))+\
103.                 sum(power(rSet[:,-1]-tree['right'],2))
104.             treeMean = (tree['left']+tree['right'])/2.0
105.             errorMerge = sum(power(testData[:,-1]-treeMean,2))
106.             if errorMerge < errorNoMerge:
107.                 print "Merging"
108.                 return treeMean
109.             else:
110.                 return tree
111.         else:
112.             return tree
113. def treeForeCast(tree, inData, modelEval=regTreeEval):
114.     if not isTree(tree):
115.         return modelEval(tree, inData)
116.     if inData[tree['spInd']]>tree['spVal']:
117.         if isTree(tree['left']):
118.             return treeForeCast(tree['left'], inData, modelEval)
119.         else:
120.             return modelEval(tree['left'], inData)
121.     else:
122.         if isTree(tree['right']):
123.             return treeForeCast(tree['right'], inData, modelEval)
124.         else:
125.             return modelEval(tree['right'], inData)
126. def createForeCast(tree, testData, modelEval=regTreeEval):
127.     m=len(testData)
128.     yHat = mat(zeros((m,1)))
129.     for i in range(m):
```

```
130.         yHat[i,0]=treeForeCast(tree, mat(testData[i]), modelEval)
131.     return yHat
132.     """
133. myData2 = loadDataSet(r"ex2.txt")
134. myMat2 = mat(myData2)
135. tree2 = createTree(myMat2, ops=(0,1))
136. print tree2
137. myData2Test = loadDataSet(r"ex2test.txt")
138. myMat2Test = mat(myData2Test)
139. print prune(tree2, myMat2Test)
140. '''
141. trainMat = mat(loadDataSet('bikeSpeedVsIq_train.txt'))
142. testMat = mat(loadDataSet('bikeSpeedVsIq_test.txt'))
143. myregTree=createTree(trainMat, ops=(1,20))
144. mymodTree=createTree(trainMat, modelLeaf, modelErr, (1,20))
145. yregHat=createForeCast(myregTree, testMat[:,0])
146. ymodHat=createForeCast(mymodTree, testMat[:,0], modelTreeEval)
147. regCo = corrcoeff(yregHat, testMat[:,1], rowvar=0)[0,1]
148. modCo = corrcoeff(ymodHat, testMat[:,1], rowvar=0)[0,1]
149. print "reg", regCo
150. print "model", modCo
```

上一篇 [机器学习经典算法7-线性回归](#)
下一篇 [机器学习经典算法9-k-means](#)

顶 踩
0 0

我的同类文章

统计分析与数据挖掘（16）

- [R的基本使用\(1\)](#)
- [关于机器学习课程的小记](#)
- [机器学习经典算法3-朴素贝叶斯](#)
- [机器学习经典算法4-logistic回归](#)
- [机器学习经典算法7-线性回归](#)

- [利用matlab进行简单的贝叶斯网络构建](#)
- [机器学习经典算法2-决策树](#)
- [机器学习经典算法5-支持向量机SVM](#)
- [机器学习经典算法6-AdaBoost](#)

[更多](#)

主题推荐 机器学习 算法

猜你在找

- 有趣的算法（数据结构）

机器学习经典算法详解及Python实现——线性回归
- 数据结构和算法

机器学习经典算法8-树回归
- Winform数据库编程:ADO.NET入门

机器学习经典算法详解及Python实现——CART分类决
- 《C语言/C++学习指南》加解密密篇（安全相关算法）

机器学习经典算法7-线性回归
- Java经典算法讲解

机器学习经典算法详解及Python实现——Logistic回



小米笔记本



app外包



小米笔记本电脑



手游回合制游戏



超薄笔记本

查看评论

3楼 [柳枫大人](#) 2015-11-18 15:48发表



第52行应该是有问题的，dataSet里面的元素是matrix类型的，无法用set

Re: [rururur](#) 2015-11-26 11:17发表



回复柳枫大人：那你是怎么解决的？我也遇到同样的问题了

2楼 [ourarewe](#) 2015-11-09 22:44发表



请问楼主用的是什么数据，那个ex00之类又是什么

1楼 [upc2whu](#) 2014-03-16 21:36发表



请问楼主主程序都运行了吗？为什么我运行有错误。

```
def binSplitDataSet(dataSet, feature, value):  
    mat0 = dataSet[nonzero(dataSet[:, feature]>value)[0],:]  
    mat1 = dataSet[nonzero(dataSet[:, feature]<=value)[0],:]  
    return mat0, mat1
```

这一段代码中，mat0 和 mat1对应的feature列是不是应该去掉？否则会无限递归下去

Re: [rururur](#) 2015-11-26 11:20发表



回复upc2whu：我感觉这块也有错，他目的是想把dataSet分成两部分

Re: [abnormal_zzb](#) 2014-12-17 19:43发表



回复upc2whu：你运行是什么错误？

您还没有登录,请[\[登录\]](#)或[\[注册\]](#)

* 以上用户言论只代表其个人观点，不代表CSDN网站的观点或立场

核心技术类目

全部主题 Hadoop AWS 移动游戏 Java Android iOS Swift 智能硬件 Docker
OpenStack VPN Spark ERP IE10 Eclipse CRM JavaScript 数据库 Ubuntu NFC
WAP jQuery BI HTML5 Spring Apache .NET API HTML SDK IIS Fedora XML
LBS Unity Splashtop UML components Windows Mobile Rails QEMU KDE Cassandra
CloudStack FTC coremail OPhone CouchBase 云计算 iOS6 Rackspace Web App
SpringSide Maemo Compuware 大数据 aptech Perl Tornado Ruby Hibernate ThinkPHP
HBase Pure Solr Angular Cloud Foundry Redis Scala Django Bootstrap

[公司简介](#) | [招贤纳士](#) | [广告服务](#) | [银行汇款帐号](#) | [联系方式](#) | [版权声明](#) | [法律顾问](#) | [问题报告](#) | [合作伙伴](#) | [论坛反馈](#)

网站客服 杂志客服 微博客服 webmaster@csdn.net 400-600-2320 | 北京创新乐知信息技术有限公司 版权所有 | 江苏乐知网络技术有限公司 提供商务支持
京 ICP 证 09002463 号 | Copyright © 1999-2014, CSDN.NET, All Rights Reserved

