下载 GitChat 论坛







登录 注册

k近邻算法 (KNN) 及kd树简介 (KD-Tree)

2016年04月24日 17:40:19 标签: k近邻算法 / KNN / kd树 / kd-tree / 机器学习

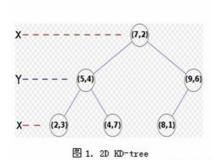
5380

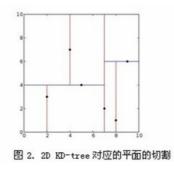
在使用k近邻法进行分类时,对新的实例,根据其k个最近邻的训练实例的类别,通过**多数表决**的方式进行预测。由于 意味看整体模型变得复杂,容易<mark>发生过拟合</mark>,即如果邻近的实例点恰巧是噪声,预测就会出错,极端的情况是k=1,称为 那么无论输入实例是什么,都简单地预测它属于训练集中最多的类,这样的模型过于简单。经验是,k值一般去一个比较小

大时尤其重要。k近邻法的最简单实现是线性扫描,这时要计算输入实例与每一个训练实例的距离,当训练集很大时,计算 😘 附,这种方法是不可行的。为了提高k近邻搜索的效率,**可以考虑使用特殊的结构存储训练数据**,以减少计算距离的 次数。具体方法有很多,这里介绍kd树方法。

二人例

先以一个简单直观的实例来介绍k-d树算法。假设有6个二维数据点 $\{(2,3),(5,4),(9,6),(4,7),(8,1),$ (7,2) }, 数据点位于二维空间内(如图2中黑点所示)。k-d树算法就是要确定图2中这些分割空间的分割线(多维空间即 为分割平面,一般为超平面)。下面就要通过一步步展示k-d树是如何确定这些分割线的。





k-d树算法可以分为**两大部分**,一部分是有关k-d树本身这种数据结构**建立**的算法,另一部分是在建立的k-d树上如何进 行最邻近查找的算法。

2.构造kd树

kd树是一种对k维空间中的实例点进行存储以便对其进行快速搜索的树形数据结构。kd树是二叉树,表示对k维空间的 一个划分。构造kd树相当于不断地用垂直于坐标轴的超平面将k维空间进行切分,构成一系列的k维超矩形区域。kd树的每 一个节点对应于一个k维超矩形区域。k-d树是一个二叉树,每个节点表示一个空间范围。下表给出的是k-d树每个节点中主 要包含的数据结构。

域名	数据类型	描述
Node-data	数据矢量	数据集中某个数据点,是n维矢量(这里也就是k维)
Range	空间矢量	该节点所代表的空间范围
split	整数	垂直于分割超平面的方向轴序号
Left	k-d树	由位于该节点分割超平面左子空间内所有数据点所构成的k-d树
Right	k-d树	由位于该节点分割超平面右子空间内所有数据点所构成的k-d树
parent	k-d树	父节点

从上面对k-d树节点的数据类型的描述可以看出构建k-d树是一个逐级展开的递归过程。下面给出的是构建k-d树的伪 码。

[cnn]



suibianshen2012

原创 粉丝 喜欢 39 13

访问量: 427 等级: 博客 5 积分: 4626 排名: 7640



达内可靠吗









他的最新文章

python向上取整, 向下取整

Linux 文件末尾包含^M字符

N-gram详解分析

Python+Hadoop Streaming实现 educe (如何给map和reduce的脏 参数)

pyhton列表合并、追加操作exten

文章分类

Python

Eclipse

jquery

maven

spring security

MyEclipse

展开~

文章存档

2017年10月

2017年9月

2017年8月

2017年7月

2017年6月

2017年5月

.

展开~

登录

淮册

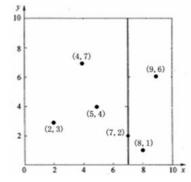
```
输出: Kd, 类型为k-d tree
      1.If Data-set为空,则返回空的k-d tree
      2.调用节点生成程序:
         (1) 确定split域:对于所有描述子数据(特征矢量),统计它们在每个维上的数据方差。假设每条数据记录为64维,
      可计算64个方差。<mark>挑选出最大值</mark>,对应的维就是split域的值。数据方差大表明沿该坐标轴方向上的数据分散得比较开,在
      这个方向上进行数据分割有较好的分辨率;
         (2) 确定Node-data域:数据点集Data-set按其第split域的值排序。位于正中间的那个数据点被选为Node-data。此
      时新的Data-set' = Data-set \ Node-data(除去其中Node-data这一点)。
      3.dataleft = {d属于Data-set' && d[split] ≤ Node-data[split]}
6
        Left_Range = {Range && dataleft}
       dataright = {d属于Data-set' && d[split] > Node-data[split]}
        Right Range = {Range && dataright}
      4.left = 由(dataleft, Left_Range)建立的k-d tree, 即递归调用createKDTree(dataleft, Left_Range)。并设
\equiv
      置left的parent域为Kd:
        right = 由 (dataright, Right_Range) 建立的k-d tree, 即调用createKDTree (dataleft, Left_Range)。并设
      置right的parent域为Kd。
```

 \odot

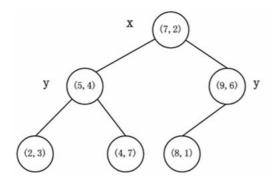
以上述举的实例来看,过程如下:

由于此例简单,数据维度只有2维,所以可以简单地给x,y两个方向轴编号为0,1,也即split={0,1}。

- (1) 确定split域的首先该取的值。分别计算x,y方向上数据的方差得知x方向上的方差最大,所以split域值首先取 0,也就是x轴方向;
- (2) 确定Node-data的域值。根据x轴方向的值2,5,9,4,8,7排序选出中值为7,所以Node-data = (7,2)。这样,该节点的分割超平面就是通过(7,2)并垂直于split = 0(x轴)的直线x = 7;
- (3) 确定左子空间和右子空间。分割超平面x=7将整个空间分为两部分,如下图所示。x<=7的部分为左子空间,包含3个节点{(2,3), (5,4), (4,7)}; 另一部分为右子空间,包含2个节点{(9,6), (8,1)}。



如算法所述,k-d树的构建是一个递归的过程。然后对左子空间和右子空间内的数据重复根节点的过程就可以得到下一级子节点(5.4)和(9.6)(也就是左右子空间的'根'节点),同时将空间和数据集进一步细分。如此反复直到空间中只包含一个数据点,如下图所示。最后生成的k-d树如下图所示。



3.搜索kd树

在k-d树中进行数据的查找也是特征匹配的重要环节,其目的是检索在k-d树中与查询点距离最近的数据点。这里先以一个简单的实例来描述最邻近查找的基本思路。

星号表示要查询的点(2.1,3.1)。**通过二叉搜索,顺着搜索路径很快就能找到最邻近的近似点**,也就是叶子节点(2,3)。而找到的叶子节点并不一定就是最邻近的,**最邻近**肯定距离查询点更近,**应该位于以查询点为圆心且通过叶子节点的圆域内。**为了找到真正的最近邻,还需要进行'回溯'操作:算法沿搜索路径反向查找是否有距离查询点更近的数据点。此例中先从(7.2)点开始讲行二叉查找。然后到达(5.4),最后到达(2.3),此时搜索路径中的节点为小干(7.2)和(5.

Python3字典中items()和python2ritems()有什么区别

□ 27964

Python如何输出两位小数的百分数

24963

python中arange()和linspace()区 23395

Apriori算法简介---关联规则的频复 算法

19050

关于Python中两个列表的比较

15191

python2.x和python3.x中raw_inpnput()区别

14453

linux 安装问题make: 没有指明目 找不到makefile。 停止

14189

python中range()、xrange()和npe()区别

10322

linux下安装teamviewer方法

383

ireport解决中文乱码问题

2 8269

人工智能课程



联系我们



请扫描二维码联系

webmaster@c 400-660-0108

■ QQ客服
■ 客

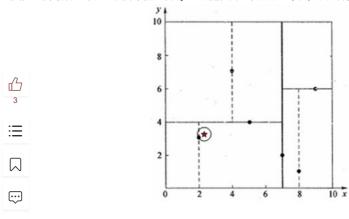
经营性网站备案信息

网络110报警服务

中国互联网举报中心

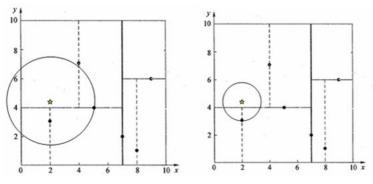
北京互联网违法和不良信息举报中心

(5,4) ,并判断在该父节点的其他子节点空间中是否有距离查询点更近的数据点。以(2.1,3.1)为圆心,以0.1414为半径画圆,如下图所示。发现该圆并不和超平面y = 4交割,因此不用进入(5,4)节点右子空间中去搜索。



再回溯到 (7,2) ,以 (2.1,3.1) 为圆心,以0.1414为半径的圆更不会与x = 7超平面交割,因此不用进入 (7,2) 右子空间进行查找。至此,搜索路径中的节点已经全部回溯完,结束整个搜索,返回最近邻点 (2,3) ,最近距离为0.1414。

一个复杂点了例子如查找点为(2, 4.5)。同样先进行二叉查找,先从(7,2)查找到(5,4)节点,在进行查找时是由y=4为分割超平面的,由于查找点为y值为4.5,因此进入右子空间查找到(4,7),形成搜索路径<(7,2),(5,4),(4,7)>,取(4,7)为当前最近邻点,计算其与目标查找点的距离为3.202。然后回溯到(5,4),计算其与查找点之间的距离为3.041。以(2, 4.5)为圆心,以3.041为半径作圆,如下图左所示。可见该圆和y=4超平面交割,所以需要进入(5,4)左子空间进行查找。此时需将(2,3)节点加入搜索路径中得<(7,2),(2,3)>。回溯至(2,3)叶子节点,(2,3)距离(2,4.5)比(5,4)要近,所以最近邻点更新为(2,3),最近距离更新为1.5。回溯至(7,2),以(2,4.5)为圆心1.5为半径作圆,并不和x=7分割超平面交割,如下图右所示。至此,搜索路径回溯完。返回最近邻点(2,3),最近距离1.5。



k-d树查询算法的伪代码如下所示。

```
[cpp]
算法: k-d树最邻近查找
输入: Kd,
          //k-d tree类型
    target //查询数据点
输出: nearest, //最邻近数据点
            //最邻近数据点和查询点间的距离
   dist
1. If Kd为NULL,则设dist为infinite并返回
2. //进行二叉查找,生成搜索路径
                                //Kd-point中保存k-d tree根节点地址
  Kd_point = &Kd;
  nearest = Kd_point -> Node-data; //初始化最近邻点
  while (Kd point)
     push (Kd_point) 到search_path中; //search_path是一个堆栈结构,存储着搜索路径节点指针
/*** \  \, \textbf{If Dist (nearest, target)} \  \, \rightarrow \, \textbf{Dist (Kd\_point -> Node-data, target)}
        nearest = Kd_point -> Node-data; //更新最近邻点
        Max_dist = Dist(Kd_point, target); //更新最近邻点与查询点间的距离 ***/
     s = Kd point -> split;
                                            //确定待分割的方向
     If target[s] <= Kd_point -> Node-data[s]
                                            //进行二叉查找
        Kd_point = Kd_point -> left;
        Kd point = Kd point ->right:
  nearest = search_path中最后一个叶子节点; //注意:二叉搜索时不比计算选择搜索路径中的最邻近点,这部分已被
  Max_dist = Dist (nearest, target);
                                    //直接取最后叶子节点作为回溯前的初始最近邻点
```

```
//确定分割方向
            s = back_point -> split;
           If Dist (target[s], back_point -> Node-data[s]) < Max_dist //判断还需进入的子空间
               If target[s] <= back_point -> Node-data[s]
                  Kd_point = back_point -> right; //如果target位于左子空间,就应进入右子空间
                  Kd_point = back_point -> left; //如果target位于右子空间,就应进入左子空间
               将Kd_point压入search_path堆栈;
           If Dist (nearest, target) > Dist (Kd_Point -> Node-data, target)
                                                           //更新最近邻点
               nearest = Kd_point -> Node-data;
4
               Min_dist = Dist(Kd_point -> Node-data, target); //更新最近邻点与查询点间的距离
```

i维数较大时,直接利用k-d树快速检索的性能急剧下降。假设数据集的维数为D,一般来说要求数据的规模N满足条 远大于2的D次方,才能达到高效的搜索。

···

转目: http://blog.csdn.net/qll125596718/article/details/8426458

Q 目前您尚未登录,请 登录 或 注册 后进行评论



olzhy 2017-12-30 23:45

回复 1楼

关于kd树 看看这篇文章吧 说的非常详细 还有代码实现 https://leileiluoluo.com/posts/kdtree-algorithm-and-i mplementation.html

Kd-Tree算法原理简析



⚠ u012423865 2017年08月22日 23:31 🕮 1075

本文讲解旨在讲解Kd-Tree(K-demension tree)的一些粗浅的原理,以及其在计算机视觉的一些应用,既是总结了自己,也 是分享给大家,希望有所帮助。 Kd-Tree是从BST(Bin...

[学习笔记]kd-tree



kd-tree学习笔记

Python凭什么又火了?!

看看现在的新闻......凭什么又火了你心里还没数吗?



KD-Tree入门



※ Zeyu King 2015年06月20日 07:39

□ 1541

经过这几天研究kd-tree,我可以说kd-tree就是按照基本的思路随便写就可以了吗?以二维平面为例,在二维平面上有若干 点,我们如何建立kd-tree?第一层以x坐标的中位数将所有点分为两部分,...

详解KDTree



卿 qp120291570 2014年11月26日 16:51 🚇 19647

简介kd树(k-dimensional树的简称),是一种分割k维数据空间的数据结构。主要应用于多维空间关键数据的搜索(如:范 围搜索和最近邻搜索)。一个KDTree的例子上图的树就是一棵KDTree,形...

最近邻查找算法kd-tree



🙀 pipisorry 2016年08月12日 10:12 🕮 13127

http://blog.csdn.net/pipisorry/article/details/52186307海量数据最近邻查找的kd-tree简介本文的主要目的是讲一下如何 创建k-d tree对特征...

加入CSDN,享受更精准的内容推荐,与500万程序员共同成长!

登录

恭喜: 一个公式教你秒懂天下英语

老司机教你一个数学公式秒懂天下英语



k ree: k近邻查询和范围查询



想象 下我们有如下两个任务: 我现在想骑一辆小黄车,我想查找离我最近的k辆小黄车. 找到百度地图中显示在屏幕上区域中 的所有酒店 这两个任务均可以用kd-tree来解决 kd-tree 主要两个用途....



k _□ ee小结



来省队集训被吊打,于是无聊学了一下kd-tree,挺好玩的东西。 kd-tree是一种支持查询平面最近点(其实是k维,但这种题

Kd-tree的用法



qq 33690156 2016年09月06日 19:46 🚇 375

Kd-Tree算法原理和开源实现代码 本文介绍一种用于高维空间中的快速最近邻和近似最近邻查找技术——Kd-Tree (Kd 树)。Kd-Tree,即K-dimensional tree,是...

kD-tree 的C语言实现 带有史上最全的注释和解释

原文转自CSDN http://blog.csdn.net/daringpig/article/details/7652891 语言cstructtreefloatinsert ...



😭 sherry_0009 2013年06月23日 09:32 🕮 1047

kdtree,kd-tree

2008年04月15日 18:54 465KB 下载



kD-tree 的C语言实现 带有史上最全的注释和解释

kdtree的原理就是基于二叉树的形式,将高维空间用超矩形进行划 😭 daringpig 2012年06月11日 16:13 🖫 17134 分.其主要用途是用来求解高维空间中最近邻的值. 下面是kdtree. h文件, 是kdtree数据结构的头文件 #i...

计算机学校排名

上海计算机专业学校学校排名

百度广告



K相邻和Kd tree



(anbing510 2013年01月28日 15:18 🕮 6593

动机 先前写了一篇文章《SIFT算法研究》讲了讲SIFT特征具体是如何检测和描述的,其中也提到了SIFT常见的一个用途就是 物体识别,物体识别的过程如下图所示: 如上图(a),我们先...

kd-tree实现KNN



■ so_sunshine 2016年09月21日 21:48 ■ 2211

转载: http://www.cnblogs.com/v-July-v/archive/2012/11/20/3125419.html 从K近邻算法、距离度量谈到KD树、...

OpenCV——KD Tree (介绍完整的flann邻近搜索)

关于OpenCV的FLANN库中KD-Tree搜索的文章并不多。 FLANN搜索的流程包括索引的建立和搜索。而所谓kd-tree搜索, 指的是在建立索引这一步骤中建立的是kd-tree索引。 本文将结合...



KD-Tree 总结



🥮 u010697167 2014年03月18日 17:53 🕮 1959

hdu 2966 In case of failure 求每个点的最短距离点的距离。直接建KD树,然后查询最近点距离就可以了。此题是比着点 击打开链接写的,可以用另一种不用删除和插入的写法。 /...

加入CSDN,享受更精准的内容推荐,与500万程序员共同成长!

登录

kd树简介 在matlab下VLFeat中的kd-tree使用 🧣 damant 2016年03月17日 21:51 🖽 4884



先简要介绍knn——K近邻算法和kd-tree——kd树,然后介绍matlab

环境中有关使用kd树的函数。k-d树(k-dimensional树的简称),是一种分割k维数据空间的数据结构。主要应用于多...

程序员不会英语怎么行?



ku-tree in Python

ki 🔍 索改进 in Python

lipengcn 2016年03月20日 20:33 🕮 2501

ka-tree c语言代码



/* This file is part of ``kdtree'', a library for working with kd-trees. Copyright (C) 2007-2009 Joh...

从K近邻算法、距离度量谈到KD树、SIFT+BBF算法

从K近邻算法、距离度量谈到KD树、SIFT+BBF算法前言 前两日,在微博上说:"到今天为止,我至少亏欠了3篇文章待 写: 1、KD树; 2、神经网络; 3、编程艺术第28章。你看到, blog内的文章与...



🧶 v_JULY_v 2012年11月20日 16:31 🖫 159332

KNN (三) --KD树详解及KD树最近邻算法

k-d树搜索最近点,在opencv中使用FLANN算法,,其包 含: 1: 建树 2.查询程序见下: #include "kdtree.h" #incl ude #include #includ...

→ App_12062011 2016年07月21日 21:24 □ 14359