

一种优化初始中心点的 K-means 算法^{*}

汪 中 刘贵全 陈恩红

¹(中国科学技术大学 计算机科学与技术系 合肥 230027)

²(安徽省计算与通讯软件重点实验室 合肥 230027)

摘 要 针对 K-means 算法所存在的问题, 提出一种优化初始中心点的算法. 采用密度敏感的相似性度量来计算对象的密度, 启发式地生成样本初始中心. 然后设计一种评价函数——均衡化函数, 并以均衡化函数为准则自动生成聚类数目. 与传统算法相比, 本文算法可得到较高质量的初始中心和较稳定的聚类结果. 实验结果表明该算法的有效性和可行性.

关键词 K-means 算法, 密度, 初始中心点, 均衡化函数
中图法分类号 TP 311

A K-means Algorithm Based on Optimized Initial Center Points

WANG Zhong, LIU Gui-Quan, CHEN En-Hong

¹ (Department of Computer Science and Technology, University of Science and Technology of China, Hefei 230027)

² (Key Laboratory of Software in Computing and Communications of Anhui Province, Hefei 230027)

ABSTRACT

Aiming at the problems of K-means algorithm, a method is proposed to optimize the initial center points through computing the density of objects. Thus, the initial center of the samples can be built in a heuristic way. Then, a new evaluation function is proposed, namely equalization function, and consequently the cluster number is generated automatically. Compared with the traditional algorithms, the proposed algorithm can get initial centers with higher quality and steadier cluster results. Experimental results show the effectiveness and feasibility of the proposed algorithm.

Key Words K-means Algorithm, Density, Initial Center Point, Equalization Function

1 引 言

聚类分析是一种无指导的机器学习方法, 是数

据挖掘中一个非常活跃的分支, 具有广泛地应用. 它基于“物以类聚”的原理, 把一组个体按照相似性归成若干类别, 使得属于同一个类别的个体之间的差

^{*} 国家自然科学基金资助项目 (No. 60775037)

收稿日期: 2007-09-30 修回日期: 2008-01-07

作者简介 汪中, 男, 1984年生, 硕士研究生, 主要研究方向为数据挖掘、机器学习. E-mail: wzspk@mail.usc.edu.cn 刘贵全, 男, 1970年生, 博士, 副教授, 主要研究方向为数据挖掘、人工智能、网络安全等. 陈恩红, 男, 1968年生, 教授, 博士生导师, 主要研究方向为数据挖掘、机器学习、网络信息处理等.

别尽可能的小,而不同类别上的个体间的差别尽可能的大^[1].为了对数据对象进行聚类,目前已有大量经典的算法涌现,Han等人归纳了基于划分、层次、密度、网络和模型的五大聚类算法^[2].

K-means算法是目前应用最为广泛的一种基于划分的聚类方法,具有简洁和快速的优点.尤其对于数值属性的数据,它能较好地体现聚类在几何和统计学上的意义.但是原始的K-means算法也存在一些缺陷:1)算法要求用户事先给定k值,在实际中,由于缺乏经验,k值一般是难以确定的;2)对初始聚类中心敏感,对于不同的初始中心,可能会导致不同的聚类结果;3)评价函数的选取有很多.文献[3]总结常用的有效评价函数:分离系数,分离熵,紧致与分离效果函数等.但是这些函数对于求解最优的聚类数目也不是很理想.

为了克服原始K-means算法的不足,不同学者从不同角度提出一系列K-means算法的变体.Huang提出一种基于K-means的变量自动加权聚类算法,使得聚类问题中的变量选择得到改进^[4].Dhillon等人则通过调整迭代过程中重新计算聚类中心的方法使其性能得到提高^[5].Zhang等人利用权值对数据点进行软分配以调整其迭代优化过程^[6].针对高维数据,杨风召^[7]等人提出一个新的相似性度量函数,较好地克服了传统的距离函数在高维空间的缺点.Sarafis则将遗传算法应用于K-means的目标函数构建中,并提出一个新的聚类算法^[8].文献[9]采用基于密度的方法,另外在微聚的时候采用一种更为有效的剪枝方法,提高聚类的精度和效率.Kaufman^[10]等人提出一种启发式的方法,通过估计数据点的局部密度作为样本初始值.文献[11]提出一种初始化K-means的谱方法.文献[12]提出一种密度敏感的相似性度量,将其引入谱聚类得到密度敏感的谱聚类算法.

针对传统K-means算法存在的问题,本文从3个方面提出改进:1)利用谱图理论的相似性思想,在密度敏感的相似度量的基础上,设计一种新的基于密度的初始化方法;2)提出一种新的评价函数——均衡化函数;3)自动生成最优聚类数目.

2 K-means算法的基本思想

设样本数据集,用 $X = \{x_i | i = 1, 2, \dots, n\}$, $C_j (j = 1, 2, \dots, k)$ 表示聚类的k个类别, $c_j (j = 1, 2, \dots, k)$ 表示初始的聚类中心,两个数据对象之间的欧式距离

$$d(x_i, x_j) = \sqrt{(x_i - x_j)^T (x_i - x_j)}$$

聚类中心

$$c_j = \frac{1}{n_j} \sum_{x_i \in C_j} x_i$$

K-means算法的核心思想是通过迭代把数据对象划分到不同的簇中,以求目标函数

$$E = \sum_{i=1}^k \sum_{j=1}^{n_j} d(x_i, c_j)$$

最小化,从而使生成的簇尽可能地紧凑和独立.算法描述如下.

算法1 K-means算法

输入 聚类的数目k和包含n个对象的数据集

输出 满足目标函数值最小的k个簇

step1 从n个数据对象中任意选择k个作为初始聚类中心.

step2 循环step3和step4直到目标函数E不再发生变化为止.

step3 根据每个聚类对象的均值,计算每个对象与这些中心对象的距离,并且根据最小距离重新对相应的对象进行划分.

step4 重新计算每个聚类的均值.

3 基于密度的初始中心点算法

由于初始中心是随机选择的,所以导致聚类结果的随机性,并且一般不会得到全局最优解.因此怎样找到一组初始中心点,从而获得一个较好的聚类效果并消除聚类结果的波动性对K-means算法具有重要意义.

在一般的基于贪心算法的初始中心点搜索过程中,由于仅基于距离因素,往往找到许多孤立点作为中心点.实际上,对于初始中心点,除了希望分布得尽量散之外,还希望这些中心点具有一定的代表性,即具有较高的密度.而现有的密度估计法^[10]是基于输入数据服从高斯混合分布的假设提出的,认为输入域的密集区域中包含自然的聚类划分,从而通过辨识输入域的密集区域得到聚类的初始中心.局限性在于仅在凸形结构的数据集上有好的聚类效果,不适合具有任意复杂形状的聚类问题.

针对高斯函数的局限性,本文提出一种基于密度的初始中心点选择算法.采用密度敏感的相似性度量函数来满足聚类的全局一致性特征,逐个从数据集中选取代表点,直至k个为止.其中第一个代表点是选择密度最大的点,其余的点根据一个启发式规则依次从剩余的代表点中选取.

3.1 基本定义

定义 1(聚类对象的密度) 已知样本事务数据库 $D = \{x_1, x_2, \dots, x_n\}$, 其中对象 x_i 的密度记作 $\text{density}(x_i)$, 即

$$\text{density}(x_i) = \sum_{j=1}^n \frac{1}{m \sum_{k=1}^{P_{ij}} \sigma^{d(x_k, x_{k+1})} - 1}.$$

采用密度敏感的相似性度量来计算密度, 其中, $d(x_i, x_j)$ 表示对象 x_i 和 x_j 的欧式距离, σ 为密度参数, P_{ij} 为连接数据点 x_i 和 x_j 之间的所有路径, 代表连接数据点 x_i 和 x_j 路径中数据的个数.

定义 2(聚类对象的邻域) 对于任意样本对象 x_i 以 x_i 为中心, R 为半径的圆形区域, 称该区域为对象 x_i 的邻域, 记为

$$\delta = \{x \mid 0 < d(x, x_i) \leq R\}.$$

定义 3 聚类对象的邻域半径

$$R = \frac{\text{aver}(D)}{n^{\text{coe}R}},$$

其中, $\text{aver}(D)$ 表示所有样本对象间距离的平均值, n 是样本对象的数目, $\text{coe}R$ 是邻域半径调节系数.

3.2 基于密度的初始化中心点算法

算法的基本思想是首先在样本事务数据库 D 中选择密度最大的点作为第一个初始中心点, 然后在数据集 D 中删去该点及其邻域内的所有对象, 再按同样的方法确定第二个初始中心点, 循环执行直到初始中心点集 M 中有 k 个点.

算法描述如下.

step1 根据定义 1 计算所有样本对象的密度 $\text{density}(x_i)$, 初始化中心点集 M 为空, 即 $M = \{\}$.

step2 选择密度最大的样本对象

$$x_{\max} = \max\{x_i \mid x_i \in D, i = 1, 2, \dots, n\}$$

作为第一个初始中心点, 加到中心点集 M 中, $M = M \cup \{x_{\max}\}$, 并从样本数据库 D 中删去该对象, 即 $D = D - \{x_{\max}\}$. 根据定义 2 和定义 3 计算 x_{\max} 邻域内的所有的样本对象, 并从样本数据库 D 中删去.

step3 重复执行 step2 直到初始中心点集中有 k 个中心点, 即 $|M| = k$

step4 输出初始中心点集 M 算法结束.

3.3 算法说明

文献 [13] 也指出, 聚类的困难在于算法要适应各种情况, 聚类应用都是根据经验和实验来确定参数. 由于密度参数 σ 和邻域半径调节系数 $\text{coe}R$ 是未知的, 所以它们的取值因实验数据的不同而不同, 并且会影响初始中心点的选择和最终的聚类结果, 所以密度参数和邻域半径系数均有经验值给出. 其中

密度参数 $\sigma > 1$, 邻域半径调节系数 $0 < \text{coe}R < 1$, 邻域半径 R 要尽可能地反应样本的空间分布, 过大过小均不能有效达到最优的聚类效果.

相对于现有的初始化中心方法, 本方法的优点如下.

1) 与文献 [10] 采用的高斯函数相比, 本文设计的密度可在距离测度上直接计算相似度. 这样可满足数据在同一流形上的数据点具有较高的相似性, 另外还可克服高斯函数对尺度敏感的缺陷.

2) 对于不同的数据集, 选择适合的经验值.

4 基于初始化中心点和均衡化评价函数的 K-means 算法

一般对聚类有效性进行评价的方法是相对准则. 其基本思想是, 在同一个数据集上, 用同一种聚类算法取不同的输入参数从而得到相应的聚类结果, 对这些不同的聚类结果, 再应用已定义的有效性函数作比较来判断最优划分. 基于数据集分类结构先验知识的评价方法除了判定聚类结果的有效性, 更主要的是检验一个聚类算法在特定数据集上的聚类效率. 本文通过对聚类评价准则的分析, 提出一种新的评价函数——均衡化评价函数, 对发现聚类结果有良好的性能, 能够找出最优的聚类个数.

4.1 均衡化评价函数

大多数为聚类设计的评价函数都着重两个方面: 每个簇的内部应该是紧凑的, 各个簇之间的距离应该尽可能地远. 实现这种概念的一种直接方法就是考察聚类 C_i 的类内差异 $w(C_i)$ 和类间差异 $b(C_i)$.

类内差异 $w(C_i)$ 可用多种距离函数来定义, 最简单的就是计算类内的每一点到它所属类中心的距离的平方和:

$$w(C_i) = \sum_{j=1}^k w(C_{ij}) = \sum_{j=1}^k \sum_{x \in C_{ij}} d(x, C_{ij})^2.$$

类间差异 $b(C_i)$ 定义为不同聚类中心间的距离:

$$b(C_i) = \sum_{i \leq j \leq k} d(C_i, C_j)^2.$$

类内差异衡量聚类的紧凑性, 类间差异衡量不同聚类之间的距离, 聚类的总体质量可以被定义为类内差异和类间差异的组合. 基于这一思想, 我们提出均衡化评价函数.

定义均衡化评价函数为类内差异和类间差异平方和的二次方根:

$$J(c, k) = \sqrt{w(C_i) \times w(C_i) + b(C_i) \times b(C_i)}. \quad (1)$$

利用均衡化评价函数作为准则函数, 有效地均

衡类内差异和类间差异的不协调性. 当均衡化评价函数达到最小时, 将得到最优的空间聚类结果, 其中 k 的选择如下:

$$\min J(c, k), \quad k=1, \dots, K \tag{2}$$

4.2 基于初始化中心点和均衡化函数的 K-means 算法

根据经验规则^[14] 知, k 的结果不超过 \sqrt{n} 改进的算法事先不需要给定 k 算法从 $1 \sim \lfloor \sqrt{n} \rfloor$ (\sqrt{n} 取整) 循环执行 K-mean 算法, 以均衡化的评价函数作为准则函数, 搜索评价函数值最小的并记下相应的 k 值, 即最优聚类结果的数目. 所以改进算法根据评价函数最小值自动生成聚类数目, 其中选择初始中心点时, 采用基于密度的方法, 获取高质量的初始中心.

算法 2 基于初始化中心点和均衡化函数的 K-mean 算法

输入 具有 n 个对象的样本数据库及经验值 σ 和 coef

输出 最优的 k 个簇, 使均衡化函数准则最小

```
step 1 for  $i = 1$  to  $\lfloor \sqrt{n} \rfloor$  do
    step 1.1 调用基于密度的初始中心点算法来
    确定  $i$  个对象作为初始中心.
    step 1.2 计算簇中的平均值, 将每个样本对
    象赋给最近的簇.
    step 1.3 更新簇的平均值.
    step 1.4 根据式 (1) 来计算评价函数  $J(c, k)$ 
    直到其收敛为止, 否则 goto step 1.2.
step 2 根据式 (2) 来搜索均衡化函数  $J(c, k)$ 
值最小的, 记下相应的  $k$  值即最优聚类数目.
step 3 结束.
```

4.3 本文算法分析

为了将原始数据的分类特征与得到的聚类相对比来评价算法的有效性, 采用划分相似测度来评价算法. 其中划分相似测度以 Fowlke 和 Mallow 测度 (FM) 来作为评价聚类结果实验的手段, 相似测度在 $0 \sim 1$ 之间取值, 其值越大表明聚类结果越能符合聚类特点, 聚类精确度的计算结果在实验部分给出.

本文算法的时间复杂度. 首先算法在初始中心点时要花费一定的时间去选取初始中心, 主要花费在计算所有样本距离的均值, 时间复杂度为 $O(n \log n)$. 由于计算的数据量较大, 可采用存储到硬盘的方式来节省内存空间. 其次在自动生成聚类数目的时候, 由于 $k \leq \sqrt{n}$ 只需进行 \sqrt{n} 次循环, 所以算法的总共时间复杂度为 $O(n \sqrt{n} \log n)$.

算法通过 \sqrt{n} 循环, 得到算法采用基于密度的算法消除初始值对聚类结果的影响. 同时对奇异点和噪声起到一定的抑制和消除作用, 均衡化函数采用距离平方的形式, 考虑类内差异和类间差异的内在关系. 相对现有算法, 有很大的改进.

5 实验分析

为了验证本文算法的有效性, 我们进行两类实验. 其中, 第一类实验采用二维数据, 以便对算法的聚类效果进行直观分析, 同时我们也与文献 [14] 进行比较. 第二类实验采用 UC 数据库的 4 组数据, 从初始中心和评价函数两个方面进行实验, 同时与实际中心以及传统算法进行比较.

5.1 实验一及分析

表 1 给出一个样本事务数据库, 并对它实施上述改进的 K-means 算法.

表 1 样本事务数据库
Table 1 Sample transaction database

序号	属性 1	属性 2	序号	属性 1	属性 2	序号	属性 1	属性 2
1	1	1	5	3	4	9	8	3
2	1	2	6	4	5	10	9	2
3	2	2	7	5	4	11	10	3
4	3	1	8	2	5	12	9	4

该表中有 12 个样本对象, 为了表达样本数据库的空间效果, 把表中的属性对应到坐标轴上. 图 1 给出了样本的空间分布图, 其中 m 是样本数据库的中心. 当 $k=3$ 时, 首先根据基于密度的初始中心点算法, 选择 3 个点作为初始中心, 即密度最大的 3 个点 A B C 点. 然后将每个样本对象赋给最近的中心点, 一次操作就可以得到最优的 3 个类, 加快聚类的过程.

图 2 给出均衡化函数 $J(c, k)$ 与类内差异 $w(c)$ 和类间差异 $b(c)$ 的趋势图. 按照本文算法, 我们计算得当 k 为 1, 2, 3 时, 均衡化函数的值分别为 $J(c, k)$ 为 39.59 20.86 15.56. 即当 $k=3$ 时, 均衡化函数值最小. 其实均衡化函数由类内差异和类间差异两部分组成, 均衡化评价函数 $J(c, k)$ 的最小值对应 $k=3$ 也正是样本数据库的最优聚类数目, 其中类内差异 $w(c)$ 随着 k 值的增大, 呈下降趋势. 而类间差异 $b(c)$ 随着 k 值的增大呈上升趋势, 并且当 k 等于样本个数时, 类间差异 $b(c)$ 达到最大.

相比之下,在文献[14]中,提出距离评价函数 $s(c,k) = b(c) + w(c)$ 的概念. 给出当类内差异 $w(c)$ 和类间差异 $b(c)$ 相等时, 聚类结果最优. 应用在本例中, 对应的 k 值并不是整数值, 所以有点偏差, 但是本文提出的均衡化评价函数 $J(c,k)$ 可适用于其它样本数据库, 所以要比距离评价函数 $s(c,k)$ 优越. 图 2 同时给出两种评价函数的比较图, 其中 $s(c,k)$ 为距离评价函数. 可以看出, 类内差异 $w(c)$ 和类间差异 $b(c)$ 值相等时, k 的取值偏向 4 均衡化评价函数 $J(c,k)$ 的结果反应了类内差异和类间差异的相互依赖关系.

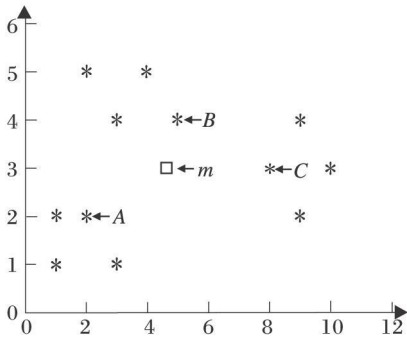


图 1 样本空间分布图
Fig 1 Spatial distribution of samples

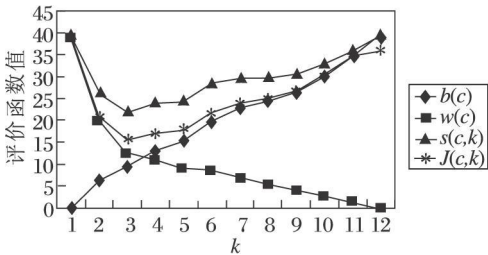


图 2 评价函数趋势及比较
Fig 2 Trends and comparison of evaluation function

5.2 实验二及分析

为了验证本文算法的有效性, 选用实验平台 Windows XP, Matlab 7.0 语言编程环境, Intel Pentium 1.40GHz CPU, 256MB 内存, 采用 UC 数据库上的 Iris, Balance scale, New-thyroid, Haberman 这 4 组数据作为测试数据. UC 数据库是一个专门用于测试机器学习、数据挖掘算法的数据库, 库中的数据都有确定的分类, 因此可以直观表示聚类结果的质量.

5.2.1 初始中心实验结果

为了说明初始中心选取的有效性, 采用 Iris 数据, 它的实际聚类中心位置分别为

(5.00 3.42 1.46 0.24), (5.93 2.77 4.26 1.32), (6.58 2.97 5.55 2.02), 其中对于参数 $\alpha \in \mathbb{R}$ 根据经验值, 当 $\alpha = 0.13$ 时, 聚类效果较好. 将本文算法所得的实验结果与 Iris 数据实际中心进行, 结果见表 2.

从表 2 中可以看出, 本文算法的实验结果与实际中心非常接近, 并且误差平方和较小, 说明该算法对 Iris 数据比较有效.

表 2 本文算法实验结果与原始的 Iris 数据中心比较
Table 2 Comparison of results between proposed algorithm and original data centers in Iris

本文算法的初始化结果	Iris 数据实际中心	误差平方和
(5.007 3.416 1.461 0.238)	(5.00 3.42 1.46 0.24)	
(5.895 2.728 4.276 1.352)	(5.93 2.77 4.26 1.32)	0.0568
(7.025 3.140 5.826 2.180)	(6.58 2.97 5.55 2.02)	

5.2.2 均衡化函数实验结果

使用以上 4 组测试数据对传统算法和本文算法分别测试 5 次, 本文算法取其平均值, 得到结果如表 3. 由于传统的 K-means 算法初始中心点是随机选取的, 从而稳定性较差, 导致聚类质量也不高. 而本文算法由于采用基于密度初始化中心点算法, 根据对象的密度分布找出聚类中心, 从而找到对象分布密集的区域, 并且采用新的评价函数, 实验表明本文算法明显优于传统算法.

表 3 传统算法和本文算法的性能比较
Table 3 Performance comparison between traditional algorithm and proposed algorithm

				%
算法	Iris	Balance scale	New-thyroid	Haberman
传统算法 1	83.67	50.48	76.07	53.96
传统算法 2	67.33	46.96	68.19	56.00
传统算法 3	76.00	42.88	81.25	52.00
传统算法 4	86.00	47.24	64.78	54.88
传统算法 5	58.33	43.88	71.04	51.00
本文算法	88.00	52.88	86.95	80.08

6 结束语

传统的 K-means 算法要求用户事先给定 k 值, 限制了很多实际应用. 初始中心点随机选择, 导致聚类结果的不稳定性, 并且容易陷入局部极值点. 常用

的评价函数对于求解最优的聚类数目也不是很理想. 本文结合谱图理论的思想, 提出一种基于密度敏感的相似度量方法来初始样本值, 有效确定 k 个代表性强的对象作为初始中心, 加快聚类过程. 并且采用新的评价函数. 实验结果证明, 该算法大大提高聚类的质量和稳定性.

参 考 文 献

- [1] Mao Guojun, Duan Lijuan, Wang Shi, et al. Principle and Algorithm of Data Mining. Beijing, China: Tsinghua University Press, 2005 (in Chinese)
(毛国君, 段立娟, 王实, 等. 数据挖掘原理与算法. 北京: 清华大学出版社, 2005)
- [2] Han J, Kamber M. Data Mining Concepts and Techniques. Orlando, USA: Morgan Kaufmann Publishers, 2001
- [3] Shi Zhongzhi. Knowledge Discovery. Beijing, China: Tsinghua University Press, 2004 (in Chinese)
(史忠植. 知识发现. 北京: 清华大学出版社, 2004)
- [4] Huang J Z, NG M K, Rong Hongqiang, et al. Automated Variable Weighting in Kmeans Type Clustering. IEEE Trans on Pattern Analysis and Machine Intelligence, 2005, 27(5): 657—668
- [5] Dhillon I S, Guan Yuqiang, Kogan J. Refining Clusters in High Dimensional Text Data // Proc of the 2nd SIAM Workshop on Clustering High Dimensional Data. Arlington, USA, 2002, 59—66
- [6] Zhang B. Generalized K-Harmonic Means. Dynamic Weighting of Data in Unsupervised Learning // Proc of the 1st SIAM International Conference on Data Mining. Chicago, USA, 2001, 1—13
- [7] Yang Fengzhao, Zhu Yangyong. An Efficient Method for Similarity Search on Quantitative Transaction Data. Journal of Computer Research and Development, 2004, 41(2): 361—368 (in Chinese)
(杨凤召, 朱扬勇. 一种有效的量化交易数据相似性搜索方法. 计算机研究与发展, 2004, 41(2): 361—368)
- [8] Sarafis J, Zalazala A M S, Trinder P W. A Genetic Rule-Based Data Clustering Toolkit // Proc of the Congress on Evolutionary Computation. Honolulu, USA, 2002, 1238—1243
- [9] Ma J, Perkins S. Time Series Novelty Detection Using One Class Support Vector Machines // Proc of the International Joint Conference on Neural Networks. Portland, USA, 2003, III, 1741—1745
- [10] Kaufman L, Rousseeuw P J. Finding Groups in Data: An Introduction to Cluster Analysis. New York, USA: John Wiley & Sons, 1990
- [11] Qian Xian, Huang Xuanjing, Wu Lide. A Spectral Method of Kmeans Initialization. Acta Automatica Sinica, 2007, 33(4): 342—346 (in Chinese)
(钱线, 黄萱菁, 吴立德. 初始化 Kmeans 的谱方法. 自动化学报, 2007, 33(4): 342—346)
- [12] Wang Ling, Bo Liefeng, Jiao Licheng. Density Sensitive Spectral Clustering. Acta Electronica Sinica, 2007, 35(8): 1577—1581 (in Chinese)
(王玲, 薄列峰, 焦李成. 密度敏感的谱聚类. 电子学报, 2007, 35(8): 1577—1581)
- [13] Rui Xu, Wunsch D I J. Survey of Clustering Algorithms. IEEE Trans on Neural Networks, 2005, 16(3): 645—678
- [14] Li Yongsen, Yang Shanlin, Ma Xijun, et al. Optimization Study on K-Value of Spatial Clustering. Journal of System Simulation, 2006, 18(3): 573—576 (in Chinese)
(李永森, 杨善林, 马溪骏, 等. 空间聚类算法中的 K 值优化问题研究. 系统仿真学报, 2006, 18(3): 573—576)