

文章编号: 1007 4619(2006)05 0715-07

基于均值-标准差的 K 均值初始聚类中心选取算法

张文君^{1,2,3}, 顾行发^{1,2,3}, 陈良富^{2,3}, 余涛^{2,3}, 许华^{2,3}

(1 电子科技大学 自动化工程学院, 四川 成都 610054 2 遥感科学国家重点实验室, 中国科学院 遥感应用研究所, 北京 100101
3 国家航天局 航天遥感论证中心, 北京 100101)

摘要: 遥感图像分类是遥感图像处理中长期存在的一个难题, 针对不同的传感器图像, 不同的应用需求, 选择合适的分类算法非常重要。在分类中不仅要考虑分类的精度, 而且要考虑分类效率。本文研究了 K 均值算法的初始聚类中心的选择对算法本身聚类精度及效率的影响, 提出了一种高效高精度的初始聚类中心选取方案, 实验结果表明, 利用该算法进行地表分类, 效率比 ENVI 的 K -Means(K 均值)模块高。

关键词: 遥感图像分类; K 均值; 聚类

中图分类号: TP751.1 文献标识码: A

An Algorithm for Initializing of K -Means Clustering Based on Mean-standard Deviation

ZHANG Wen-jun^{1,2,3}, GU Xing-fa^{1,2,3}, CHEN Liang-fu^{2,3}, YU Tao^{2,3}, XU Hua^{2,3}

(1. School of Automation Engineering, the University of Electronic Science and Technology of China, Shidian Chengdu 610054, China
2. State Key Laboratory of Remote Sensing Science, Jointly Sponsored by the Institute of Remote Sensing Applications of Chinese Academy of Sciences and Beijing Normal University, Beijing 100101, China
3. The Center for National Spaceborne Demonstration, Beijing 100101, China)

Abstract Remote Sensing Image Classification is always a difficult problem. It's very important to use proper algorithms for different images. Not only precision but also efficiency must be considered. By researching in the relations between the initial means of clusters and the efficiency of clustering, we proposed a method of computing the initial means of clusters for K -Means Clustering. It is proved to be more efficient than the K -Means Clustering module of ENVI when it's used in Remote Sensing Classification.

Key words remote sensing image classification; clustering; K -Means

1 引 言

遥感技术是人类研究地球资源信息的重要手段。随着资源环境信息的爆炸性增长, 地理信息系统 (GIS) 在土地资源管理和环境监测等方面发挥着

越来越重要的作用, 当前, GIS 正在与遥感 (RS)、全球定位系统 (GPS) 紧密结合, 向 3S 一体化方面发展。遥感分类是遥感地理信息系统中的关键技术之一。针对不同的遥感影像不同的需求, 需要利用不同的分类算法来达到高效高精度的遥感分类。遥感分类分为监督分类与非监督分类, 至今已有很多种

收稿日期: 2006-04-10 修订日期: 2006-05-26

基金项目: 国家重大基金项目 (编号: 2002CB412506), 中国科学院百人计划项目 (编号: KZCX0415), 国家教育部留学回国人员科研启动基金重点项目 (编号: HX040013), 国防科学技术工业委员会项目 (编号: KJSX0401)。

作者简介: 张文君 (1982—), 女, 电子科技大学在读硕士生。主要从事遥感图像处理方面的研究。

成熟的算法,如非监督分类方法有动态聚类算法, K 均值法,ISODATA等;监督分类方法有最小距离法,最大似然比法;到20世纪80年代后期,人工神经网络的各种模型开始诞生,发展到今天,已经广泛运用于遥感图像非监督分类中。

在遥感图像分类的应用领域里,由于经常会出现对分类图像所代表的区域不了解,即没有先验知识的情况,因此,非监督分类显得尤为重要,其中,聚类分析在非监督分类中应用非常广泛。

聚类分析就是将数据对象分组成为多个类或簇,在同一个簇中的对象之间具有较高的相似度,而不同簇中的对象差别较大。通过聚类,人们能够识别密集的和稀疏的区域,从而发现数据的整体分布模式,还能找到数据间的有趣的相互关系。关于聚类分析目前已经有 K 均值,CURE,ISODATA等很多算法,而且在实践中得到了应用。 K 均值算法由MaQueen^[1,2]首先提出,是解决聚类问题的一种经典算法。该算法具有简单、快速并且能够有效地处理大量数据的优点。但是, K 均值算法的聚类质量完全依赖于初始解的选择,它的执行结果与数据的输入次序有关^[3]。这里针对 K 均值算法做出改进,对该算法初始聚类中心值的选取做了研究,考虑了参与聚类数据的均值、标准差与数据分布的关系,由均值-标准差决定初始聚类中心,即初始解。文章的主要思想是根据两个标准:分类结果的质量和该算法对初始中心选取的敏感度来评价初始值选取的优劣,找到最优解。实验结果表明,该算法在精度跟收敛速度上,都比其他随机选取初始中心的方法高。

2 K 均值算法

2.1 算法描述

K 均值方法是基于划分的聚类方法。它在目前的聚类分析中应用最为广泛。其基本思想为:对于给定的聚类数目 K ,首先随机创建一个初始划分,然后采用迭代方法通过将聚类中心不断移动来尝试着改进划分。

具体描述如下:输入一个数据集和一个整数 K (簇的个数),输出的一个划分 $p_k = \{C_1, \dots, C_k\}$ 。

定义^[4] K 均值聚类问题:假设 N 个数据集合 $x = \{x_1, \dots, x_N\}$,是待聚类数据,其中 $x_j = \{x_{j1}, \dots, x_{jd}\} \in R^d$, $j = 1, \dots, N$ 。 K 均值聚类问题是要找到 X 的一个划分 $p_k = \{C_1, \dots, C_k\}$,使目标函数 $f(p_k) =$

$\sum_{i=1}^k \sum_{x \in C_i} d(x, m_i)$ 最小。其中, $m_i = \frac{1}{n_i} \sum_{x \in C_i} x_i$ 表示第 i 个簇的中心位置, $i = 1, \dots, k$, n_i 是簇 C_i 中数据项的个数, $d(x, m_i)$ 表示 x_i 到 m_i 的距离。

通常的 K 均值的处理流程如下:随机地将数据集划分成 k 个簇 $C(1, \dots, k)$,计算每个簇的平均值 $m_i, i = 1, \dots, k$ 作为簇中心,然后将每个数据项按照其与各个簇中心的距离,重新分配到最近的簇。再计算每个簇的中心。若中心位置发生改变,则重复这个过程,否则算法停止。找到一个局部极小的划分。

K 均值算法的具体流程如下:

输入: N 个 d 维待分类数据 $\{x_1, x_2, \dots, x_n\}$,其中 $x_i = \{x_{i1}, \dots, x_{id}\}$,待分类的簇数 K ;

输出: K 个簇,使得所有数据与离其最近的簇中心相异度总和最小;

步骤 1:随机选择 K 个初始聚类中心 $\{c_1, c_2, \dots, c_k\}$,其中, $c_j = \{c_{j1}, \dots, c_{jd}\}$;选择聚类最大迭代次数 l ;确定迭代结束的最大收敛系数 T ;

步骤 2:根据欧氏距离公式,计算每个数据到各簇的距离,将各数据分到具有最小距离的簇中,其中距离计算公式为:

$$d(x_i, m_j) = \sqrt{\sum_{l=1}^d (x_{il} - m_{jl})^2}, \quad i = 1, \dots, N; \quad j = 1, \dots, k;$$

$d(x_i, m_j)$ 为第 i 个矢量数据到第 j 个聚类的距离。

步骤 3:重新计算 K 个聚类的中心值 $\{m_1, m_2, \dots, m_k\}$,计算公式为:

$$m_j = \frac{1}{n_{x_i \in C_j}} \sum_{x_i \in C_j} x_i, \quad m_j = \{m_{j1}, \dots, m_{jd}\}, \quad l = 1, \dots, d;$$

m_j 为第 j 个聚类的聚类中心。

步骤 4:检验聚类操作是否应该结束。

若迭代次数等于 l 则结束聚类,否则计算该次迭代的各个聚类收敛距离,若每个簇的收敛距离都小于给定的参数 T ,则结束;否则,继续迭代,迭代次数加一,转向步骤 2。收敛距离的计算公式为:

$$t_j(k) = \sqrt{\sum_{l=1}^d (m_{jl}(k) - m_{jl}(k-1))^2}$$

式中, $j = 1, \dots, K$; $l = 1, \dots, d$; k 为当前迭代次数。

2.2 算法存在的问题

尽管该算法已经得到了广泛应用,但它也存在着一些不可避免的问题。最重要的几点如下:

- 与很多聚类方法一样, K 均值算法是在假设参与聚类的数据中聚类数 K 知道的前提下进行的,不一定与实际的类别数一致。

· 在迭代技术上, K 均值算法对初始聚类条件特别敏感(初始聚类和数据输入秩序)。

· K 均值算法可能造成局部最小解。算法使初始聚类中心到最终的聚类结果有唯一的一一对应的映射关系。

为了克服输入参数值 K 与真实值不同的缺陷, 可以选取一种普遍的方法: 在聚类时输入几个 K 值多次聚类, 最后分析每种聚类结果, 选取较优的输出。关于初始聚类中心与最后结果具有确定性对应的问题几乎是所有聚类算法的普遍问题, 如果能找到最优初始聚类中心, 则可在一定程度上提高聚类的质量, 这个问题也可以忽略了。因此, 围绕着聚类的质量和收敛的敏感性, 重点讨论初始聚类中心的选取。

3 初始条件选取与比较

3.1 随机选取初始聚类中心

在上述算法分析中, 初始条件决定了聚类收敛速度和聚类的精度, 如果随机选取, 则分类结果精度和效率也存在着随机性。实验表明, 选择不同的初始聚类中心, 所得的聚类离散度以及收敛速度都是

不一样的。

为了验证, 对于同一组数据, 选取不同的聚类初始中心值对收敛速度和聚类精度的影响, 采用随机数生成方法, 生成下列 4 组数据, 将这 4 组数据分别执行 K 均值算法。针对遥感图像的特定需求, 将另外 3 组数据的值控制在 0—255 之间, 生成数据的统计值见表 1。

表 1 测试数据统计特征值

Table 1 Test data's statistic information

数据编号	数据个数	中心值	标准差
1	1000	992.996	407.915
2	10000	125.844	73.8687
3	100000	126.379	73.8115
4	1000000	126.807	73.5421

对上面的 4 组数据分别用 K 均值算法做 4 次聚类实验, 其中每次聚类选取的聚类中心都是随机选取, 设定聚类数为 5 聚类的收敛系数为 0 即各聚类中心位置始终不再改变时结束该操作, 算法收敛时所需要的迭代次数和花费时间见表 2。

表 2 聚类结果

Table 2 Results of clustering

测试号	数据编号	初始类中心	聚类结果中心	迭代次数	时间/(ms)
F-1	1	{176 395 882 958 1581}	{289 101 616 691 923 772 1224 94 1604 95}	22	15
F-2	1	{548 580 804 856 949}	{289 101 616 691 923 772 1224 94 1604 95}	25	16
F-3	1	{308 637 1083 1419 1717}	{361 497 722 21 1009 62 1313 71 1660 9}	15	16
F-4	1	{35 48 94 197 296}	{289 101 616 691 923 772 1224 94 1604 95}	37	15
2-1	2	{26 51 61 135 203}	{24 2948 75 5333 126 438 176 792 228 825}	31	288
2-2	2	{19 29 74 85 221}	{24 2948 75 5333 126 438 176 792 228 825}	34	187
2-3	2	{23 30 205 203 245}	{26 3525 80 4062 131 708 180 496 230 129}	19	125
2-4	2	{41 44 107 214 235}	{24 2948 75 5333 126 438 177 149 229 177}	10	63
3-1	3	{45 146 159 200 228}	{25 8847 78 4711 129 763 180 442 230 169}	27	1641
3-2	3	{8 52 71 118 213}	{23 3924 72 0375 121 868 174 039 227 708}	27	1578
3-3	3	{18 19 38 129 151}	{23 3924 72 0375 121 868 174 039 227 708}	32	1890
3-4	3	{144 157 179 246 248}	{25 8847 78 4711 129 763 180 442 230 169}	34	2094
4-1	4	{6 60 75 91 114}	{22 9816 71 0169 120 954 172 96 227 029}	26	18375
4-2	4	{67 89 122 204 236}	{26 4967 79 5288 130 946 181 482 230 526}	13	8906
4-3	4	{44 65 99 135 157}	{22 9816 71 0169 120 954 172 96 227 029}	19	13734
4-4	4	{38 69 73 142 185}	{22 9816 71 0169 120 465 172 467 227 029}	20	14484

由以上数据可以看出, 选取不同的初始聚类中心, 所需要的迭代循环次数不同, 从而导致系统时间开销不一样。当小数据量时, 如数据为 1000 至 10000 条, 所需的系统时间开销差距不大, 初始聚类中心的选取对时间开销并无太大的影响; 当参与聚类的数据量增大到 100000 甚至 1000000 条时, 系统时间开销的差距增大到几百甚至几千 μs 。此时, 初始聚类中心的选取严重地影响着系统时间开销。在遥感图像分类中, 一般的遥感图像像素个数都在百万以上, 并且多数情况下分类的维数为多维, 对于高光谱图像, 参与分类的维数更多, 此时初始聚类中心的选取严重地决定着分类所需要的时间。因此, 在 K 均值算法用于遥感图像分类中, 随机选取初始聚类中心值对于分类时间的影响是随机的, 必须找出一种有效的初始聚类中心, 以提高非监督分类的效率。

3.2 基于均值-标准差选取初始点

由 K 均值算法可知, 如果所选取的初始聚类中心选取在几个分布密集区域的中心, 其周围的点越容易分到最近的点, 聚类收敛越快, 所需要迭代的次数越少。其中涉及最优初始聚类中心点的选取。以二维矢量数据为例, 图 1 为该矢量数据集的分布情况, 数据主要分布在 A, B 点周围, 若要分为两类, 则选取 A, B 两点为初始聚类中心时, 其周围的点到该初始中心的距离差距越大, 越容易收敛。相反, 若选择 C, D 两点为初始中心, 收敛速度慢。

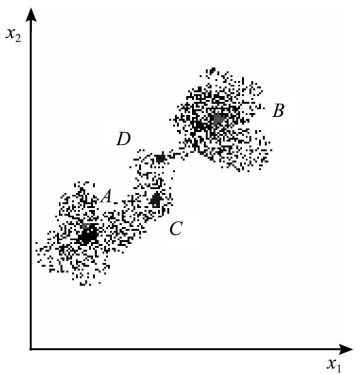


图 1 二维数据在特征空间中的分布

Fig. 1 Distribution of two dimensional data

若要分析所有数据的分布情况, 计算其分布密度, 那是非常复杂的事。根据随机函数的分布知识, 聚类的数据应主要分布在所有数据的均值附近。标准差是评价数据分布的又一重要指标, 假设所有数据的均值为 μ , 标准差为 σ , 则数据应该主要分布在

$(\mu - \sigma, \mu + \sigma)$ 之间。假设分类数为 N , 选择初始分类点为 $(\mu - \sigma, \mu + \sigma)$ 之间的 N 个等分点进行分类。设第 i 类的初始分类中心为 m_i , 则:

$$m_i = (\mu - \sigma) + \frac{2\sigma_i}{N}, i = 1, \dots, N;$$

如果参与分类的是多维数据, 如 d 维, 则每个聚类初始聚类中心的各个向量为 $(\mu_l - \sigma_l, \mu_l + \sigma_l)$ 之间, 设第 i 类聚类初始中心值为 $\{m_{i1}, m_{i2}, \dots, m_{id}\}$, 则有: $m_{il} = (\mu_l - \sigma_l) + \frac{2\sigma_l^i}{N}, i = 1, \dots, N; l = 1, \dots, d$ 。

以二维向量数据为例, 其均值, 初始聚类中心点的关系如图 2 所示。

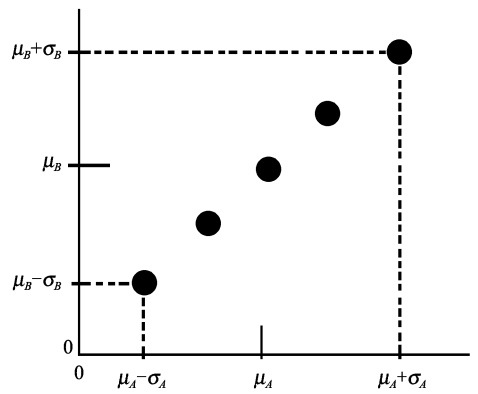


图 2 二维数据初始聚类中心

Fig. 2 5 K Means initial means of clusters

采用上述方法, 对前面的 4 组数据分类, 分类数为 5 收敛距离为 0 其分类统计结果见表 3。

3.3 比较

比较表 3 和表 2 利用均值方差选用初始聚类中心的方法进行 K 均值分类时, 算法迭代次数最少, 系统花销时间明显减少, 当参加聚类的数据达到上百万条时, 该方法所花费的时间比前述方法少 10s 以上。在时间效率上提高的同时, 聚类的精度并没有减少, 以第四组数据为例 (第四组有上百万数据, 值在 0—255 之间), 前面四次测试与最后一次考虑初始聚类中心的方法的测试结果精度及时间比较见表 4。

表 4 描述了选择不同的初始聚类中心, 设定收敛系数相同的条件下, 各聚类最终的聚合情况和时间耗费。可以看出, 每聚类的平均聚合程度相同, 而考虑用均值-方差计算初始聚类中心的方法聚类结果各类的聚合程度更均一, 而时间上的花费少很多。

表 3 聚类时间比较
Table 3 Comparing of clustering efficiency

测试号	数据编号	初始类中心	聚类结果中心	迭代次数	时间 / (ms)
1-5	1	{585.08, 789.038, 992.996, 1196.95, 1400.91}	{356.772, 705.856, 986.642, 1286.3, 1643.97}	11	16
2-5	2	{51.9757, 88.9, 125.844, 162.779, 199.713}	{25.7927, 78.9937, 129.639, 178.319, 229.177}	9	62
3-5	3	{52.5677, 89.4734, 126.379, 163.285, 200.191}	{24.8865, 75.4699, 125.736, 176.903, 228.642}	9	578
4-5	4	{53.2648, 90.0358, 126.807, 163.578, 200.349}	{24.9761, 75.9835, 126.416, 176.965, 228.528}	10	7000

表 4 聚类精度比较
Table 4 Comparing of clustering precision

测试号	聚类结果个类标准差	时间 / (ms)
4-1	{15.8597, 15.2883, 14.7221, 14.1323, 13.5653}	18375
4-2	{15.0202, 14.4321, 15.5964, 13.8332, 14.73}	8906
4-3	{14.732, 15.0202, 13.8332, 14.4321, 15.5964}	13734
4-4	{13.5653, 14.1323, 15.8597, 14.7221, 15.2883}	14484
4-5	{14.7178, 14.7139, 14.4545, 14.7213, 14.9917}	7000

4 K 均值应用于遥感图像分类

大量的数据验证表明, 将均值-方差计算初始聚类中心的 K 均值方法用在遥感影像分类中将大大提高分类速度。判断一幅遥感影像分类算法的好坏除了时间上的快慢外, 还必须考虑分类的精度问题, 即分类评估。对于无监督分类, 在不知道图像先验

知识的情况下, 判断分类精度的指标是各聚类中心间的距离和聚类内部标准差。类间距离反映了个类间的离散程度, 类内标准差反映了各类内部的聚合程度, 类间距离越大, 类内标准差越小, 则说明所得的分类结果越好, 各类聚合度越高。

下面对北京地区和苏州地区的 CBERS 卫星 IRS 传感器第 1 波段影像数据进行分类, 分为 5 类: 植被, 水体, 城市, 裸地以及其他像元, 收敛系数设置为 1。先后采用 ENV I 的 K 均值算法和本文的 K 均值算法进行分类, 影像分类的结果如图 3 所示, 左图为北京的影像图, 右图为苏州影像图。

从视觉效果上来看, 考虑初始值选取后的分类图的分类效果与 ENV I 的相当。为了比较两种算法的分类效果, 特写算法求出两种方法分类图像在原图像中的各类中心和类内标准差, 并求出两种方法所需要的迭代次数, 统计迭代次数, 聚类结果各类中心, 聚类结果各类标准差, 结果见表 5。

表 5 分类结果效率精度比较

Table 5 Comparing of efficiency and precision

(a) 迭代次数

(a) The times of iterating

分类工具	北京地区	苏州地区
	迭代次数	迭代次数
ENV I	14	11
本文算法	11	8

(b) 分类结果各类均值

(b) The means of clusters

分类工具	北京地区各聚类中心					苏州地区各聚类中心				
	类 1	类 2	类 3	类 4	类 5	类 1	类 2	类 3	类 4	类 5
ENV I	0.2501	65.08	82.79	97.06	112.87	0	43.82	55.21	67.08	76.88
本文算法	0.2291	62.30	81.18	96.25	111.93	0	42.64	53.83	66.15	76.12

(c) 分类结果各聚类标准差

(c) The standard deviations of clusters

分类工具	北京地区各聚类标准差					苏州地区各聚类标准差				
	类 1	类 2	类 3	类 4	类 5	类 1	类 2	类 3	类 4	类 5
ENV I	2.4592	7.6187	4.4931	4.1212	8.9156	0	3.8985	3.0108	2.9996	5.2551
本文算法	2.3226	7.6742	5.0034	4.1297	8.8608	0	3.4742	3.3375	3.0382	5.2854

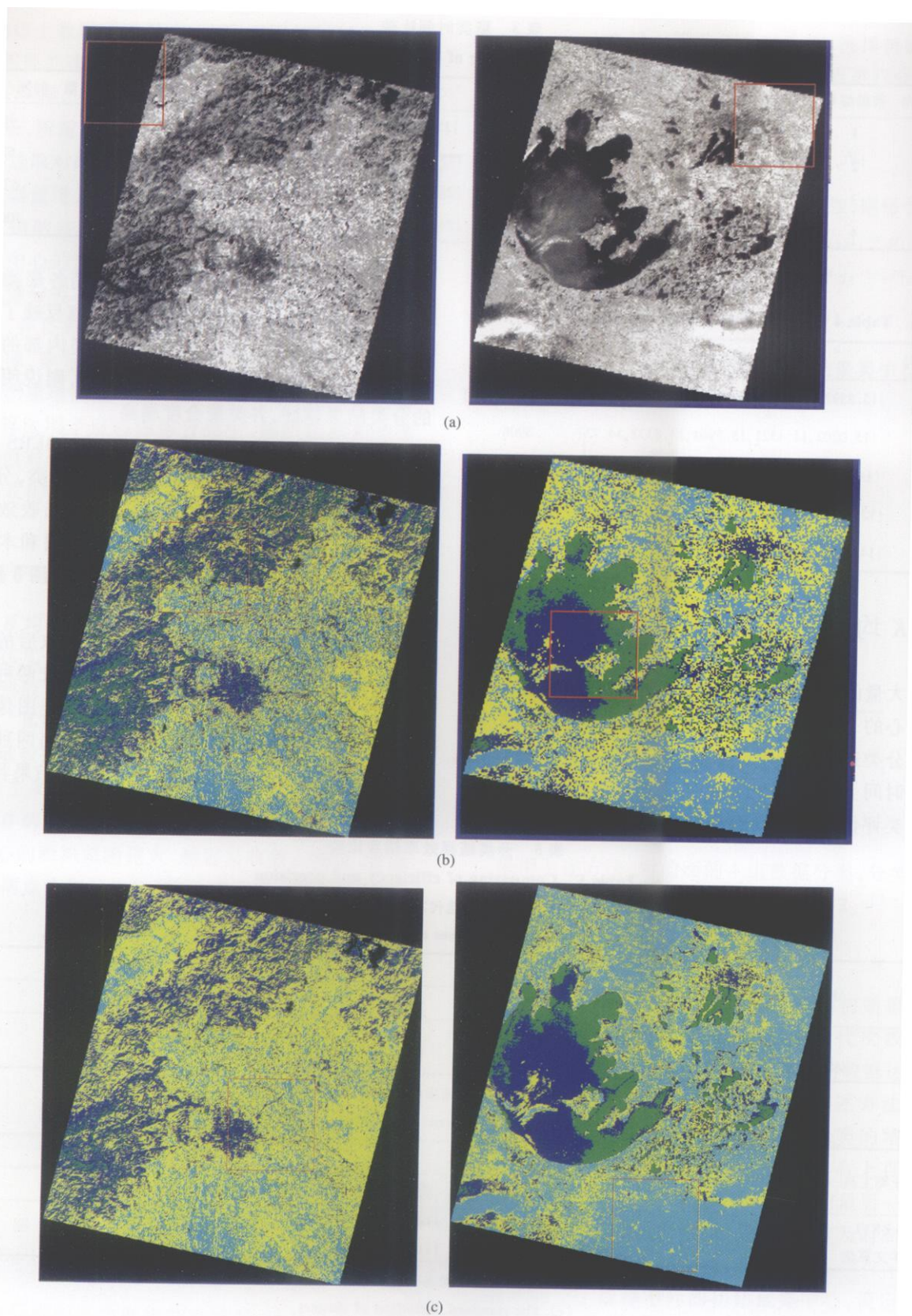


图 3 分类结果

(a) 原始图像; (b) ENV I 的 K 均值分类结果; (c) 本文算法的分类结果

Fig 3 Results of image clustering

(a) Images to be classified; (b) Images classified by ENV I sK-means algorithm; (c) Images classified by the proposed algorithm

注: 苏州的影像图有云 (绿色部分)

表 5 的数据说明, 考虑用均值-方差计算初始聚类中心点的 K 均值方法在遥感图像非监督分类中既能提高分类速度, 又能保证分类的精度。

5 结束语

本文提出了基于均值-标准差的方法确定 K 均值算法聚类初始中心的思想, 经过多组随机数据及遥感图像的测试, 验证了该方法在保证聚类精度的同时, 聚类的时间效率有很大的提高。当用于上百万至千万的多维遥感影像非监督分类时, 更能明显地改善分类时间。在今后的研究中, 可考虑初始聚类中心与聚类数据分布函数的关系, 遥感图像经过预处理后的灰度值大多数为正态分布, 考虑服从其他分布的数据, 研究其关系。另外就遥感影像来说, 今后将从直方图入手研究它与最优初始聚类中心的关系。

参考文献 (References)

- [1] Wu D, Hou Y T, Zhang Y Q. Transporting Real time Video over the Internet Challenges and Approaches[J]. *Proceeding of the IEEE*, 2000, **88**(12): 1855- 1875.
- [2] Fine Granularity Scalable MPEG4 Standards[S]. ISO /IEC JTC 1 /SC 29 WG 11 ISO /IEC JTC1 /SC 29 WG 11 N3518 Beijing 2000- 07.
- [3] Pena JM, Lozano JA, Larranaga P. An Empirical Comparison of Four Initialization Methods for the K-Means Algorithm[J]. *Pattern Recognition Letters*, 1999, **20**: 1027- 1040.
- [4] Wu J L, Zhu W X. An Iterated Local Search Algorithm for K-means Clustering[J]. *Computer Engineering and Application*, 2004, (22): 37- 41. [吴景岚, 朱文兴. 基于 K 均值的迭代局部搜索聚类算法 [J]. 计算机工程与应用, 2004, (22): 37- 41.]
- [5] Sergios Theodoridis, Konstantinos Koutroubas. *Pattern Recognition*[M]. Publishing House of Electronics Industry 2004. [Sergios Theodoridis, Konstantinos Koutroubas. *Pattern Recognition*[M]. 电子工业出版社, 2004.]