

# 基于 K 均值聚类和概率松弛法的图像区域分割

周卫星, 廖欢

(华南师范大学 物理与电信工程学院, 广东 广州 510006)

**摘要:**在进行图像区域分割时,为了减少过度分割现象,可利用 K 均值算法简单、快速并且能够有效地处理大数据库的优点及概率松弛算法并行快速且考虑空间信息的优点,同时考虑灰度信息和空间信息将两种方法相结合应用于图像的区域分割。首先利用 K 均值聚类方法将图像初步分为多个类,然后,利用迭代的概率松弛法对粗分结果进行优化,对一些疑似像素进行进一步分割和目标提取。实验结果表明,该算法比较简单且具有良好的特性。

**关键词:**图像区域分割; K 均值; 概率松弛

**中图分类号:** TP391

**文献标识码:** A

**文章编号:** 1673-629X(2010)02-0068-03

## Region - Based Image Segmentation Based on K - means and Probability Relaxation

ZHOU Wei-xing, LIAO Huan

(School of Physics & Telecommunication Engineering, South  
China Normal University, Guangzhou 510006, China)

**Abstract:** In order to reduce over-segmentation, can take advantages of K-means which is simple, fast and able to deal with large database and probability relaxation. The method of combining K-means and probability relaxation is used in this paper. First apply K-means clustering method to segment the image pixels into different regiments. Then an iterative probability relaxation operation is applied in order to optimize the coarse segmentation to further segment the uncertain pixels according to their statistic properties. Experimental results indicate that proposed method is effective for image segmentation and object extraction.

**Key words:** region-based image segmentation; K-means; probability relaxation

## 0 引言

在对图像的研究和应用中,通常需要将所关心的目标从图像背景中提取出来,图像分割就是将图像分成各具特性的区域并提取出感兴趣的区域的技术和过程。在研究图像分割的过程中,已经产生了阈值分割、边缘检测、统计学分割等多种方法<sup>[1]</sup>。但这些方法只是将图像简单地分为两类,这往往不符合实际情况。而目前有人提出的基于模糊 C 均值和概率松弛的区域分割算法又存在以下缺点:其一,在许多实际应用中,所涉及的数据大多是海量的,这必然会导致用 FCM 计算非常复杂;其二,由于聚类中心初值的设定是随机的,会导致算法在局部最小值时收敛,影响聚类的准确度。针对上述情况,考虑到 K 均值算法具有简单、快速并且能够有效地处理大数据库的优点,且分类

结果基本上不受初始聚类中心的影响,笔者提出先利用 K 均值聚类对图像像素进行初步的划分,在充分了解目标灰度信息的基础上,再利用概率松弛算法对分类结果进行更进一步的处理,从而提取出图像中真正的目标,实现对图像的成功分割。

## 1 图像分割定义

图像分割将图像中具有特殊含义的不同区域分离开来,这些区域是互不相交的,每一个区域都满足特定区域的一致性。比如对同一物体的图像,一般需要将图像中属于该物体的像素(或物体的特征像素点)从背景中分割出来,将属于不同物体的像素点离开。

分割出来的区域应该同时满足:

(1)分割出来的图像区域的均匀性和连通性。其中,均匀性指的是该区域中所有像素点都满足基于灰度、纹理、彩色等特征的某种相似性准则,连通性是指该区域内存在连接任意两点的路径。

(2)相邻分割区域之间针对选定的某种差异显著

收稿日期:2009-06-24;修回日期:2009-09-06

基金项目:广东省攻关项目(2008B080701053)

作者简介:周卫星(1958-),男,副教授,研究方向为图像处理、模式识别。

性。

(3)分割区域边界应该规整,同时保证边缘的空间定位精度。

借助集合概念对图像分割可给出如下比较正式的定义<sup>[2]</sup>:

令集合  $R$  代表整个图像区域,对  $R$  的分割可看作将  $R$  分成  $N$  个满足以下五个条件的非空子集(子区域)  $R_1, R_2, \dots, R_N$ :

- ①  $\bigcup_{i=1}^N R_i = R$ ;
- ② 对所有的  $i$  和  $j, i \neq j$ , 有  $R_i \cap R_j = \emptyset$ ;
- ③ 对  $i = 1, 2, \dots, N$ , 有  $P(R_i) = \text{TRUE}$ ;
- ④ 对  $i \neq j$ , 有  $P(R_i \cup R_j) = \text{FALSE}$ ;
- ⑤ 对  $i = 1, 2, \dots, N, R_i$  是连通的区域。

其中  $P(R_i)$  是对所有在集合  $R_i$  中元素的逻辑谓词,  $\emptyset$  代表空集。

下面先对上述各个条件分别给予简略解释:条件①指出在对一幅图像的分割结果中全部子区域的总和(并集)应能包括图像中所有像素(就是原图像),或者说分割应将图像中的每个像素都分进某个子区域中。条件②指出在分割结果中各个子区域是互不重叠的,或者说在分割结果中一个像素不能同时属于两个区域。条件③指出在分割结果中每个子区域都有独特的特性,或者说属于同一个区域中的像素应该具有某些相同特性。条件④指出在分割结果中,不同的子区域具有不同的特性,没有公共元素,或者说属于不同区域的像素应该具有一些不同的特性。条件⑤要求分割结果中同一个子区域内的像素应该是连通的,即同一个子区域内的任意两个像素在该子区域内互相连通,或者说分割得到的区域是一个连通组元。最后需要指出,实际应用中图像分割不仅要把一幅图像分成满足上面五个条件的各具特性的区域而且需要把其中感兴趣的目标区域提取出来。只有这样才算真正完成了图像分割的任务。

## 2 粗分类算法选择

模糊C均值(FCM)是解决聚类问题的一种经典算法,是基于对目标函数优化基础上的一种数据聚类方法。假设有一幅图像,它的  $n$  个像素构成模糊集  $X = (x_1, x_2, \dots, x_n)$ , 若将  $n$  个像素分成  $c$  类,则构成  $c$  个模糊子集,每个模糊子集都有一个聚类中心。对于模糊C均值,定义目标函数为:

$$J(U, V) = \sum_{k=1}^n \sum_{i=1}^c (u_{ik})^m (d_{ik})^2$$

式中,  $m \in [1, +\infty]$  为加权指数;  $u_{ik}$  表示第  $k$  个像素对第  $i$  个类的隶属度;  $d_{ik}$  为第  $k$  个像素到第  $i$  个类中

心的距离,  $U$  为模糊分类矩阵,  $V$  为聚类中心集合。

●FCM算法步骤可描述如下:

(1)确定分类数  $c$  与加权指数  $m$ ,  $m$  一般取 2, 设定迭代停止阈值  $\epsilon$  为一小正数, 初始化迭代次数  $l = 0$  和模糊分类矩阵  $U^{(0)}$ , 给出初始聚类中心  $V^{(0)}$ ;

$$(2) \text{ 将 } U^{(l)} \text{ 代入 } v_i = \frac{\sum_{k=1}^n (u_{ik})^m x^k}{\sum_{k=1}^n (u_{ik})^m}, (i = 1, 2, \dots, c)$$

计算聚类中心矩阵  $V^{(l)}$ ;

$$(3) \text{ 根据式 } u_{ik} = \frac{1}{\sum_{j=1}^c \left( \frac{d_{ik}}{d_{jk}} \right)^{\frac{2}{m-1}}}, \text{ 利用 } V^{(l)} \text{ 更新隶属$$

度矩阵  $U^{(l)}$ , 得到新的模糊分类矩阵  $U^{(l+1)}$ ;

(4)若  $\|U^{(l)} - U^{(l+1)}\| < \epsilon$ , 停止。否则置  $l = l + 1$ , 返回(2)。

算法中的  $\epsilon$  为收敛阈值,  $\epsilon > 0$ 。  $\epsilon$  是影响聚类精度和聚类速度的参数。数值过大时, 会导致算法过早收敛, 聚类结果不稳定, 特别在初始参数不确定的条件下, 这种现象更为普遍和明显。当阈值过小时, 则可能会导致过度计算, 既浪费时间, 又可能发生无法收敛的问题。

K均值算法由 MacQueen<sup>[3,4]</sup> 首先提出, 也是解决聚类问题的一种经典算法。其核心思想是随机选择  $K$  个对象, 每个对象代表一个簇的初始均值或中心。对剩余的每个对象, 计算其与各个簇中心的距离, 根据使距离指标的目标函数  $E = \sum_{k=1}^K \sum_{x \in Q_k^i} \|g(x) - \mu_i^{k+1}\|^2$  最小的原则将它划分到最相似的簇, 算法选择的相似性度量通常是欧几里德距离的倒数<sup>[5]</sup>, 也就是说两者的距离越小表示两者的相似性越大, 反之则相似性越小。然后重新计算各个簇的新均值, 更新簇中心。这个过程不断重复, 直到准则函数收敛。

●K均值算法描述如下:

(1)任意选  $K$  个初始类均值,  $\mu_1^1, \mu_2^1, \dots, \mu_K^1$ ;

(2)在第  $k$  次迭代中, 可根据下述准则将每个像素都赋给  $K$  类之一: 对所有  $i = 1, 2, \dots, K, j = 1, 2, \dots, K, i \neq j$ , 如果  $\|g(x) - \mu_i^k\| < \|g(x) - \mu_j^k\|$  则  $x \in Q_i^k$ , 设  $Q_i^k$  的新的类中心为  $Q_i^{k+1}$ ;

(3)对  $j = 1, 2, \dots, k$ , 根据式  $\mu_i^{k+1} = \frac{1}{N} \sum_{x \in Q_i^k} g(x)$  更新类均值;

(4)对所有  $i = 1, 2, \dots, K$ , 若  $\mu_i^{k+1} = \mu_i^k$ , 则终止。否则, 退回步骤(2)继续下一次迭代。

通过比较上述两种算法描述可得知 K 均值比 C 均值算法简单、快速, 对处理大数据集, 该算法是相对可伸缩的和高效率的, 且分类结果基本上不受初始聚

类中心的影响,故文中选择 K 均值聚类算法作为粗分类时的分割。

### 3 概率松弛迭代算法

松弛方法是图像处理和模式识别中一种有效的并行方法<sup>[6]</sup>,最早是由 Rosenfeld 等<sup>[7]</sup>提出。它的基本思想是图像中每一个像素的归属不仅应该由其本身来决定,而且应该受到它的领域像素的影响。所以松弛法既是建立在领域运算基础上的算法,又是依据精确的分类判据,以一个像素一个像素迭代进行的,因此它具有并行性质运算速度高的优点和串行性质逐步迭代自适应的能力。松弛法以像素为操作对象,借助一次次的迭代确定各像素的归属,每次迭代基于定量的相容性准则。

设要将有  $N$  个像素的原始集合  $S = \{S_1, S_2, \dots, S_N\}$  分成  $M$  个类,及类别集合  $C = \{C_1, C_2, \dots, C_M\}$ 。概率松弛算法引入一个名为置信度的量  $p_i(j)$  表示  $S_i \in C_j (1 \leq i \leq N, 1 \leq j \leq M)$  的概率,对于任何一个  $i$  均有  $0 \leq p_i(j) \leq 1$ 。首先确定一个初始值  $p_i^0(j)$ ,并引入一个相容系数  $c(i, j; k, l)$  来衡量一对分量  $S_i \in C_j$  与  $S_k \in C_l$  是否相容。当  $S_i \in C_j$  与  $S_k \in C_l$  这两个分类事件相互支持时,表示它们相容,原则上  $c(i, j; k, l) > 0$ ,而当  $S_i \in C_j$  与  $S_k \in C_l$  这两个分类事件相互冲突时,表示它们不相容,则  $c(i, j; k, l) < 0$ 。相容系数可以用来确定每次迭代中概率的改变量。一般情况下,概率增量  $q_i(j)$  可定义为:

$$q_i(j) = \frac{1}{N-1} \sum_{\substack{k=1 \\ k \neq i}}^N \left[ \sum_{l=1}^M c(i, j; k, l) p_k(l) \right]$$

其中,  $N$  为目标总个数;  $M$  为分类数;  $p_k(l)$  为  $S_k \in C_l$  的概率值;  $c(i, j; k, l)$  为事件  $S_i \in C_j$  与  $S_k \in C_l$  的相容系数;  $q_{ij}$  为  $S_i \in C_j$  事件的概率增量。

用以上方法确定  $p_i(j)$  的增量  $q_i(j)$  后,即可求取  $p_i(j)$  的新估计值。这里要求新的估计值为非负,简单的方法是取  $p_i(j)(1 + q_i(j))$ 。将  $p_i(j)$  的值规范化,得到:

$$q_i(j)^r = \frac{1}{N-1} \sum_{\substack{k=1 \\ k \neq i}}^N \left[ \sum_{l=1}^M c(i, j; k, l) p_k(l)^r \right]$$

$$p_i(j)^{r+1} = \frac{p_i(j)^r [1 + q_i(j)^r]}{\sum_{j=1}^M p_i(j)^r [1 + q_i(j)^r]}$$

从以上两式可知,在每次迭代过程中,对每个像素的灰度(或概率密度)来说是一次由其各个领域像素对它进行加权调整的过程,它从具有相同灰度性质的领域像素得到正的调整量,而从具有不同灰度性质的领域像素得到负的调整量。

### 4 图像区域分割的实现

(1) 随机初始化标记。

设图像  $f(x, y)$  中共有  $N$  个像素,由于松弛运算中并不考虑像素的位置,所以每个像素可用  $A(i)$  来表示,其中,  $i = 1, 2, \dots, N$ 。对预先图像分为的  $K$  个类,计算第  $k$  类 ( $k = 1, 2, \dots, K$ ) 的均值和方差。像素  $A(i)$  和第  $k$  类的马氏距离为:

$$d_{ik} = \frac{[\mu_k - A(i)]^2}{\delta_k^2}$$

则初始概率为:

$$p_i^0(k) = \frac{1/d_{ik}}{\sum_{k=1}^K 1/d_{ik}}$$

(2) 更新估计值。

对每对类  $k$  和类  $l$ ,相容性矩阵  $R^{[8]}$  定义为

$$R(k, l) = 1 \quad \text{若 } k = l$$

$$R(k, l) = 0 \quad \text{若 } k \neq l$$

如果用  $Q_i(k)$  表示类  $k$  对点  $i$  的相容性因素,用  $V(i)$  表示点  $i$  的领域,则

$$Q_i(k) = \frac{1}{n-1} \sum_{j \in V(i)} \sum_{l=1}^K R(k, l) p_j(l)$$

当考虑 8 点领域时,则  $n = 8$ 。在第  $r+1$  概率矢量可用下式计算:

$$p_i^{r+1}(k) = \frac{p_i^r(k) [1 + Q_i^r(k)]}{\sum_{l=1}^K p_i^r(l) [1 + Q_i^r(l)]}$$

(3) 迭代中止。

迭代松弛法是一种像素标记法,因此可设定一个百分量,当图像中超过这个百分量的像素都明确地标记时即  $p_i(k)$  大于这个值时,则认为迭代收敛,停止算法。

### 5 实验结果

该算法可以将图像分割成任意指定的  $K$  个区域,而且分割效果比较好。如图 1 右所示是该方法的运行结果。图中用不同的像素表示不同的区域,一共分割成了 3 个区域。



原图像

分割后的图像

图 1 分割前后对比图

(下转第 74 页)

$$\text{Weighted Semantic Similarity}(|AB|) = 1/2 (\text{Weighted Semantic Similarity}(AB) + \text{Weighted Semantic Similarity}(BA)) \quad (7)$$

这种算法把词语的权重区别加入算法,使得语义相似度的计算更加合理。

### 3 实验及结果分析

目前的智能答疑系统没有一个统一的评测标准,只能通过大量的数据来检验其准确率,把准备的《大学生计算机基础》课程的实例问题输入到系统中进行检测,计算相关度准确率的算法如下:

准确率 = 自动答对的题目 / 实验问题总数

实验结果如表 1 所示:

表 1 相关度算法准确率

测试的问题个数	未加权准确率	加权准确率
100	0.564	0.798
200	0.486	0.754
300	0.354	0.732

实验结果表明基于权重的相似度计算返回答案的准确率大概在 73.2%~79.8%,而未基于权重的相似度计算返回答案的准确率只有 35.4%~56.4%,基于权重的相似度计算返回答案的准备率得到明显提高。

然而,准确率并不是十分理想,原因主要集中在:

(1)分词的准确率,由于有些词语不能被分词模块正确划分,对最后的结果产生了较大的影响;

(2)对于知识库的依赖,也就是如果某些问题在知识库里没有相应的答案则无法返回正确答案。

因此接下来提高系统的准确率和智能性的主要任

务是分词技术的改进和知识库内容的丰富。

### 4 结束语

语句相关度计算是关键技术之一,相关度计算的应用和完善可以使用户得到的答案的准确率得到提升。文中主要详细介绍了基于权重的语句相似度计算的方法以及如何将该方法运用到智能答疑系统中,从而使得系统的准确率得到进一步的提高,也使得智能答疑系统的智能性得到进一步的体现。

#### 参考文献:

- [1] Roussionv D, Robles J. Self-learning Web question answering system[C]//World Wide Web conference(WWW2004). New York, US: ACM, 2004: 400-401.
- [2] 张亮,冯冲,陈肇雄,等.基于语句相似度计算的 FAQ 自动回复系统设计与实现[J]. 小型微型计算机系统, 2006(4): 720-722.
- [3] 张正兰,李珊.一个支持自然语言提问的智能答疑系统的实现[J]. 微机发展(现更名:计算机技术与发展), 2003, 13(12): 39-41.
- [4] 王荣波,池哲儒.基于词类串的汉语句子结构相似度计算方法[J]. 中文信息学报, 2005, 19(1): 21-29.
- [5] 揭春雨,刘源.论汉语自动分词方法[J]. 中文信息学报, 1989(1): 1-9.
- [6] 庞剑锋,卜东波,白硕.基于向量空间模型的文本自动分类系统的研究与实现[J]. 计算机应用研究, 2001, 18(9): 23-26.
- [7] 张小艳.中文自动答疑系统的研究与实现[J]. 微计算机信息, 2007, 12(3): 208-210.

(上接第 70 页)

### 6 结束语

提出了一种基于 K 均值聚类和概率松弛算法并将这种算法用于图像的区域分割过程中。在 Windows XP 操作系统环境下,用 VC6.0 实现了这种图像分割方法。由于在分割过程中充分考虑了像素点之间的特征,减小了过度分割的现象,使得分割后的结果更加贴近基于对象的分割。实验结果表明这种方法简单,运行速度快,分割效果也比较好。

#### 参考文献:

- [1] 罗希平,田捷,诸葛婴,等.图像分割方法综述[J]. 模式识别与人工智能, 1999, 12(3): 300-312.
- [2] 章毓晋.图像分割[M]. 北京:科学出版社, 2001.
- [3] MacQueen J. Some methods for classification and analysis of

multivariate observations[D]. Berkeley, Calif.: University of California Press, 1967.

- [4] Huang Z. Extensions to the k-means algorithm for clustering large data sets with categorical values[J]. Data Mining and Knowledge Discovery, 1998(2): 283-304.
- [5] 李苏梅,韩国强.基于 K 均值聚类算法的图像区域分割[J]. 计算机工程与应用, 2008, 44(16): 163-167.
- [6] Zucker S W. Relaxation Processes for Scene Labeling: Convergence, Speed, and Stability[J]. IEEE trans. on SMC, 1978(1): 41-48.
- [7] Rosenfeld A, Hummel R A, Zucker S W. Scene labeling by relaxation operations[J]. IEEE Trans. Syst. Man Cybern, 1976, 6: 420-453.
- [8] GARBAY C. Image Structure Representation and Processing: A Discussion of Some Segmentation Methods in Cytology[J]. IEEE Tran. on PAMI, 1986, 8(2): 140-146.