

# K-means 算法在电信 CRM 客户分类中的应用<sup>①</sup>

左国才<sup>1,2</sup> 杨金民<sup>1</sup> (1.湖南大学 软件学院 湖南 长沙 410082;

2.湖南软件职业学院 湖南 长沙 410205)

**摘要:** 面对电信市场竞争的日益加剧和信息技术的迅猛发展,电信运营商必须建立以“客户为中心”的管理模式,将客户进行分类,针对不同的客户,研究出相应的营销策略。数据挖掘中的 K-means 聚类算法能对大型数据集进行高效分类。对 K-means 算法进行改进,使其能够应用于复杂的电信客户关系管理,实现更加准确和全面的客户分类。

**关键词:** 电信 CRM; 数据挖掘; 客户分类; K-means; 聚类分析算法

## K-means Algorithm for CRM Customers in the Telecommunications Classification

ZUO Guo-Cai<sup>1,2</sup>, YANG Jin-Min<sup>1</sup>

(1. Software School, Hunan University, Changsha, 410082, China; 2. Hunan Vocational Institute of Software, Changsha 410205, China)

**Abstract:** With the sharp competition in the telecommunications market and the rapid development of information technology, telecommunications operators must establish a “customer-centric” management style to classify customers. For different customers, different marketing strategies are developed. Data Mining in the K-means clustering algorithm for large data sets can be efficiently classified. In this paper, K-means algorithm is used in complex telecom customer relationship management to achieve more accurate and comprehensive customer classification.

**Keywords:** telecom CRM; data mining; customer classification; K-means cluster analysis algorithm

## 1 引言

数据挖掘是数据库研究、开发和应用最活跃的分支科学之一,从大量数据中用非平凡的方法发现有用的知识和人们感兴趣的数据模式成了人们的一种自然需求。随着数据挖掘研究的蓬勃发展,出现很多数据挖掘的方法,其中聚簇是最基本的方法,它既可以独立地应用,也可以作为其他数据挖掘方法的前期工作。在聚簇方法中,k-means 算法是最著名和最常用的划分法之一。k-means 算法能有效地处理规模较大和高维的数据集合,但却只能聚簇数值数据,因为数值数据能用欧几里德距离测量不同数据对象之间的相异度。k-means 算法不能处理分类属性型数据。经过改

进后的 k-means 算法能够适应复杂的电信 CRM 中非数值数据的处理,实现更加准确和全面的客户分类。

## 2 电信客户关系管理现状分析

随着中国电信竞争格局的改变和通信技术的飞速发展,中国电信业的市场环境发生了根本性的变化,加入 WTO 后,我国电信市场将逐步对外开放,国内电信运营商正面对一个全新的、更加激烈的、国际国内全方位的市场竞争环境。电信运营商意识到在竞争越来越激烈的商业时代,资源占有成为企业生死成败的关键,客户才是企业生存和发展的根基,而如何改善客户服务,增强客户的满意度和忠诚度,提升客户

<sup>①</sup> 基金项目:国家自然科学基金(60703097,60703155);国家高技术研究发展计划(863)(2007CB310702)

收稿时间:2009-05-20



坐标。其中重载的方法 **distance** 一个用于计算该数据点与其它数据点之间的距离, 另一个用于计算数据点与代表点之间的距离, 以便确定代表点代表区域内的点集。另外需要说明, 表示数据点的类 **Point** 存放在一个一维数组 **PointSet** 中, 通过访问数组中的每个元素, 就可以得到每个数据点的信息。在扫描数据点时, 就是从数组中提取信息。

### 5.1.2 代表点的数据结构

代表点的数据结构:

```
class Reference{
float xm, ym, xs, ys;
int ns, cflag;
ArrayList ps;
public Reference(){
cflag=0;
}
public float distance(class Reference{
float d;
d=abs(xm-Reference.xm)+abs(ym-Reference.y
m);
return d;
}}
```

其中:  $x_m$  表示代表点的  $x$  坐标,  $y_m$  表示代表点的  $y$  坐标;  $x_s$  表示代表点代表区域内的所有点的  $x$  坐标之和,  $y_s$  表示代表点代表区域内的所有点的  $y$  坐标之和;  $n_s$  表示代表点代表区域内点的总数, 在算法中将要用这个参数来表示代表点的密度;  $p_s$  是一个集合的实例, 表示代表点代表区域内的点的集合; **cflag** 是在对代表点进行簇的划分时, 用于标记该代表点被划分在第几个簇中。其中方法 **distance** 用于对代表点进行簇的划分时计算各个代表点之间的距离, 如果两个代表点的距离小于或等于 2 倍的 **Radius** 则他们为邻接代表点, 即它们处于同一簇中。

### 5.1.3 改进后的 k-means 算法的 java 实现

#### (1) 什么是 k-means 聚类算法?

**K-means** 是最常用的聚类算法之一, 能有效地处理规模较大和高维的数据集合, 能对大型数据集进行高效分类, 把数据分成几组, 按照定义的测量标准, 同组内数据与其他组数据相比具有较强的相似性, 这就叫聚簇。

**k-means** 算法的效率比较高; 缺点是只能处理

数值型数据, 不能处理分类数据, 对例外数据非常敏感, 不能处理非凸面形状的聚簇。

**k-means** 算法接受输入量  $k$ ; 然后将  $n$  个数据对象划分为  $k$  个聚类以便使得所获得的聚类满足: 同一聚类中的对象相似度较高; 而不同聚类中的对象相似度较小。聚类相似度是利用各聚类中对象的均值所获得一个“中心对象”(引力中心)来进行计算的。

#### (2) k-means 算法的处理流程

- ① 从  $c$  个数据对象任意选择  $k$  个对象作为初始聚类中心;
- ② 循环③到④直到每个聚类不再发生变化为止;
- ③ 根据每个聚类对象的均值(中心对象), 计算每个对象与这些中心对象的距离; 并根据最小距离重新对相应对象进行划分;
- ④ 重新计算每个(有变化)聚类的均值(中心对象);

#### (3) k-means 算法的改进

**k-means** 算法是在数据挖掘领域中普遍应用的聚类算法, 它只能处理数值型数据, 而不能处理分类属性型数据。例如表示人的属性有: 姓名、性别、年龄、家庭住址等属性。而改进后的 **k-means** 算法就能够处理分类属性型数据。采用相异度来代替 **k-means** 算法中的距离, 相异度越小, 则表示距离越小。一个样本和一个聚类中心的相异度就是它们各个属性不相同的个数, 不相同则记为一, 最后计算一的总和。这个和就是某个样本到某个聚类中心的相异度。该样本属于相异度最小的聚类中心。

相异度测量: 设  $X$ 、 $Y$  是分类数据集中的两个对象, 该对象是  $m(x_1, x_2, \dots, x_m)$  维的, 则这两个对象之间的相异度为:

$$d_i(X, Y) = \sum_{j=1}^m \delta(X_j, Y_j); \text{ 其中 } \delta(x_j, y_j) = \begin{cases} 0(x_j = y_j) \\ 1(x_j \neq y_j) \end{cases},$$

该过程可以被描述为如下数学问题:

$$\text{最小化 } P(W, Q) = \sum_{i=1}^k \sum_{l=1}^n \sum_{j=1}^m w_{i,l} \delta(x_{i,j}, q_{i,j})$$

其中  $w_{i,l} \in W$ , 并且  $Q_l = [q_{i,1}, q_{i,2}, \dots, q_{i,m}] \in Q$ 。

算法实现的思想:

1) 确定聚簇数  $k$  以及阈值  $t$ ;

为了确定聚簇数  $k$ , 用一个基于距离计算相异度的聚簇方法聚簇样本数据, 下面是具体操作过程。

① 从原始数据中随机地取  $m$  个点作为样本数据放到集合  $S$  中, 定义一个初始阈值  $t = t_0$ 。

② 给所有的样本数据加上未被聚簇的标记, 并定义一个  $k = 0$ 。

③ 从样本数据集未被聚簇的点中选择一个初始点  $P$ 。{标记  $P$  属于聚簇  $C_k$ ; 从  $P$  开始递归地按照深度优先方式遍历各点,  $P' = \text{NEAR}(P, t)$ , 如果  $P'$  不是空就属于  $C_k$ , 否则退回到前一点接着进行搜索, 更新  $d, d$  是聚簇  $C_k$  中各对顶点之间的平均距离, 更新  $t = t_0 \times d$ ,  $t$  是阈值。}

④ 如果还有未聚簇的点, 使  $k=k+1$ , 重复(3)进行聚簇。

其中, 函数  $\text{NEAR}(P, t)$  是:

{寻找  $P$  的最近邻居  $P'$ , 即  $\text{dist}(P', P) \leq t$ ; 如果没有  $P'$  被发现就返回  $\text{NULL}$ , 否则返回  $P'$ 。}

2) 用样本数据得到的聚簇结果求聚簇模式  $Q$ , 用基于频率的方法求模式向量  $Q$  的过程如下:

求每一簇中在某一属性上对象数的百分数  $fr$ ,  $fr$  称为相对频率, 找出每个簇中相对频率最大的一个的属性作为  $Q$  向量的属性值, 即让  $n(k, j)$  是第  $k$  组属性  $j$  上的对象数,  $fr(A_j = ?ck, j/X) = n(k, j)/n$  最大者就是  $q_j$  属性。每个属性都这样求, 可以得到  $Q$  向量。

(4) 算法的实现

//找出最接近这个数据项的簇类

```
public Cluster findMostNearestNeighbor(){
int mindex=0; //最接近簇类在簇类集中的下标
ArrayList<Cluster>
cList=this.neighborClusters.clusterList; //簇类集
int length=cList.size(); //目标簇集
Cluster rCls=cList.get(0);
PointND clsCenter=rCls.center;
double mindistance=
this.rowdata.squareDistance(clsCenter);
if(mindistance==--1.0)
System.out.println("Error!");
//遍历簇集,对每个簇计算与该数据项的距离
if(length>1){
```

```
for(int i=1;i<length;i++){
Cluster cls=cList.get(i);
PointND center=cls.center;
//查找距离最小的
double newDistance=
this.rowdata.squareDistance(center);
if(newDistance<mindistance){
mindex=i;
mindistance=newDistance;}}
rCls=cList.get(mindex);
this.mostNearestCluster=rCls;
return rCls;}
//根据聚类结果对数据集分类
private void classify(Matrix theDataSet){
ArrayList<DataItem>
rows=theDataSet.rows;
for(int i=0;i<rows.size();i++){
DataItem rowdata=rows.get(i);
ClusterSet
neighbors=rowdata.neighborClusters;
Cluster
cluster=rowdata.mostNearestCluster;
int labelIndex=
neighbors.clusterList.indexOf(cluster);
labelIndex++;
rowdata.setLabel(""+labelIndex);}
//随机选择 N 个中心点进行聚类,直到分类精度超过界限
public void randomSelectTest2(int num) throws
IOException{ //进行聚类 Matrix
ranMatrix=new Matrix();
ranMatrix.rows=(ArrayList<DataItem>)dat
Set.rows.clone();
ranMatrix.height=dataSet.height;
ranMatrix.width=dataSet.width;
ClusterSet
ranClusters=this.initClusterSet(ranMatrix, num);
this.kMeansClustering(ranClusters,ranMatrix);
this.classify(ranMatrix);
this.saveFile();}
① 共有 336 组数据。KMEANS 用于存放进行标
```

准化处理过的数据, KMCLASS 用于存放进行聚类(K-Means)处理过的结果。

② jfreechart, 相应在的 jar 包。jcommon-1.0.12.jar, jfreechart-1.0.9.jar

③ 实现环境, 是在 myeclipse6.0GA, tomcat6.0.10

## 6 测试数据及运行结果分析

测试数据集是某电信公司的客户信息数据库, 数据量为 158200。实验的硬件环境: PC 计算机, CPU 为 PIV 2.0G, 内存为 2G; 软件环境: 操作系统为 Windows Professional 2000, 编程环境 Eclipse 3.3/ myeclipse6.0GA, tomcat6.0.10。

运行结果:

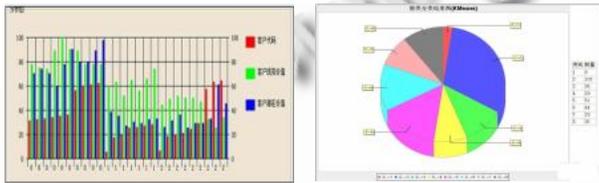


图 3 改进后的 K-means 算法运行结果

改进后的 K-means 算法在本文中主要是对客户现有价值和潜在价值进行聚类分析, 从而对客户进行分类, 制定相应的营销政策。实验结果表明, 改进后的 K-means 算法能够高效地实现客户分类, 为电信部门解决客户分类的难题。

通过以上分析, 我们最终把客户分为 8 类, 分别为: 高端商务型客户、普通商务型客户、时尚追求型客户、费用节省型客户、普通消费型客户、流动型客户、稳定型客户、潜在客户。分类的同时, 我们可以得到每个顾客被分到哪个组, 以及每组的客户选择不同产品的概率, 对于电信运营商提供个性化服务、确定客户价值、进行套餐设计和实行深度营销等都有很强的现实意义。

当然, 本研究也有一定的局限性:

① 样本量相对不大。由于数据来源问题, 样本量的选取相对不是很多, 如果有更多的样本供研究, 研究的推广价值应该更大。

② 如果能对连续很多期的数据进行研究, 应该可以对客户未来的消费情况进行预测, 也是客户流失管理所关心的内容之一。

③ 如果能结合客户的背景资料综合分析, 对现实的指导意义会更大。

综上所述, 改进后的 k-means 算法在电信客户聚类分组中的应用是相对有效的, 对于各类行业中关于客户数据的聚类分组同样具有很强的指导意义。同时希望在今后的研究中能够通过对更加完备的数据的分析, 来完善该研究方法。

## 参考文献

- Huang ZX. Extensions to the K-means algorithm for clustering large data sets with categorical values. *Data Mining and Knowledge Discovery*, 1998,2:283-304.
- Barbara D. Using self-similarity to cluster large data sets. *Data Mining and Knowledge Discovery*, 2003,7:123-152.
- Modha DS, Spangler WS. Feature Weighting. *K-Means Clustering. Machine Learning*, 2003,52:217-237.
- James D. Better together marketing research and CRM. *Marketing News*, 2002,36:1015-16.
- Hand D, Mannila H, Smyth P. 张银奎, 廖丽, 宋俊, 等译. *数据挖掘原理*. 北京:机械工业出版社, 2003.
- Michael J, Berry A, Gordon S, Linoff 袁卫等译. *数据挖掘—客户关系治理的科学与艺术*. 北京:中国财政经济出版社, 2004.
- 朱爱群. *客户关系治理与数据挖掘*. 北京:中国财政经济出版社, 2001.
- 李雄飞, 李军. *数据挖掘与知识发现*. 北京:高等教育出版社, 2005.