

聚类分析算法

聚类分析算法 (Principal component analysis, PCA) 也是一种无监督学习算法，针对没有响应变量而仅有特征变量的数据集，其主要作用就是快速分类。虽然是无监督学习算法，聚类分析也有很多应用场景，比如电商平台系统针对具有相似购买行为的用户进行聚类，针对划分好的客户类别，将某用户购买的产品在同一类别用户内相互推荐，实现精准促销；或者根据以往销售记录及其他特征对产品进行聚类，若某用户购买了一款产品，则继续向其推送同一类别内其他产品。又比如大型连锁企业依据销售情况对全辖门店进行有效分类，根据聚类结果制定生产计划、开展物流配送、实现更有效率的资源配置等。

目录

C O N T E N T S

- 1 聚类分析算法的基本原理
- 2 数据准备
- 3 划分聚类分析算法示例
- 4 层次聚类分析算法示例
- 5 习 题



PART 01

聚类分析算法的基本原理

聚类分析算法的基本原理

聚类分析是根据特征变量，按照一定的标准对样本示例进行分类。

按照分析方法的不同，聚类分析分成两个宽泛的类别，包括划分聚类分析和层次聚类分析。划分聚类分析是一种快速聚类方式，它将数据看作K维空间上的点，以距离为标准进行聚类分析，将样本分为指定的K类，包括K个平均数的聚类分析方法、K个中位数的聚类分析方法。划分聚类分析过程只限于连续数据，要求预先指定聚类数目。

层次聚类分析也称系统聚类分析，基本思路是对相近程度最高的两类进行合并，组成一个新类并不断重复此过程，直到所有的个体都归为一类，通常只限于较小的数据文件（要聚类的对象只有数百个）。

划分聚类分析

划分聚类分析方法的基本思想是将观测到的样本划分到一系列事先设定好的不重合的分组中去。划分聚类分析方法主要包括两种：一种是K个平均数的聚类分析方法，此方法的操作流程是通过迭代过程将样本示例分配到具有最接近的平均数的组，然后找出这些聚类；另一种是K个中位数的聚类分析方法，此方法的操作流程是通过迭代过程将样本示例分配到具有最接近的中位数的组，然后找出这些聚类。

以K均值聚类分析为例，K均值聚类分析的基本原理是：首先指定聚类的个数并按照一定的规则选取初始聚类中心，让个案向最近的聚类中心靠拢，形成初始分类，然后按最近距离原则不断修改不合理分类，直至合理为止。

划分聚类分析

比如用户选择 x 个特征变量参与聚类分析，最后要求聚类数为 y ，那么将由系统首先选择 y 个样本示例（当然也可由用户指定）作为初始聚类中心， x 个特征变量组成 x 维特征空间。每个样本示例在 x 维特征空间中是一个点， y 个事先选定的样本示例案就是 y 个初始聚类中心点。然后系统按照距这几个初始聚类中心距离最小的原则把样本示例分派到各类中心所在的类中去，构成第一次迭代形成的 y 类。

然后系统根据组成每一类的样本示例，计算各特征变量均值，每一类中的 x 个特征均值在 x 维特征空间中又形成 y 个点，这就是第二次迭代的聚类中心，按照这种方法依次迭代下去，直到达到指定的迭代次数或达到终止迭代的要求时，迭代停止，形成最终聚类中心。

划分聚类分析

K-均值聚类法计算量小、占用内存少并且处理速度快，因此比较适合处理大样本的聚类分析。

划分聚类分析方法与层次聚类分析方法相比在计算上相对简单且计算速度更快一些，但是它也有自己的缺点，它要求事先指定样本聚类的精确数目，这与聚类分析探索性的本质是不相适应的。

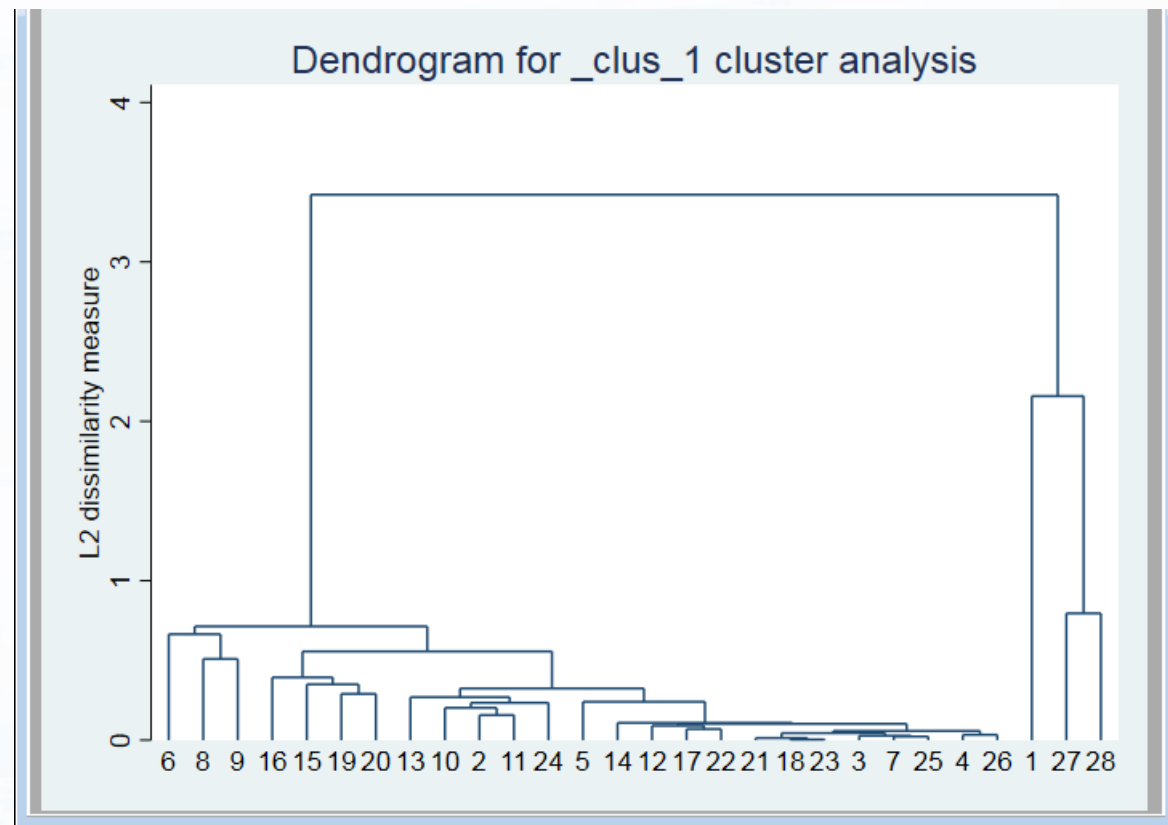
|| 层次聚类分析

层次聚类分析也称系统聚类分析，与划分聚类分析方法的原理不同，层次聚类分析的基本原理是根据选定的特征来识别相对均一的个案（变量）组，使用的算法是开始将每个个案（或变量）都视为一类，然后根据类与类之间的距离或相似程度将最近的类加以合并，再计算新类与其他类之间的相似程度，并选择最相似的加以合并，这样每合并一次就减少一类，不断继续这一过程，最终实现完全聚类，即把所有的观测样本汇集到一个组中。

|| 层次聚类分析

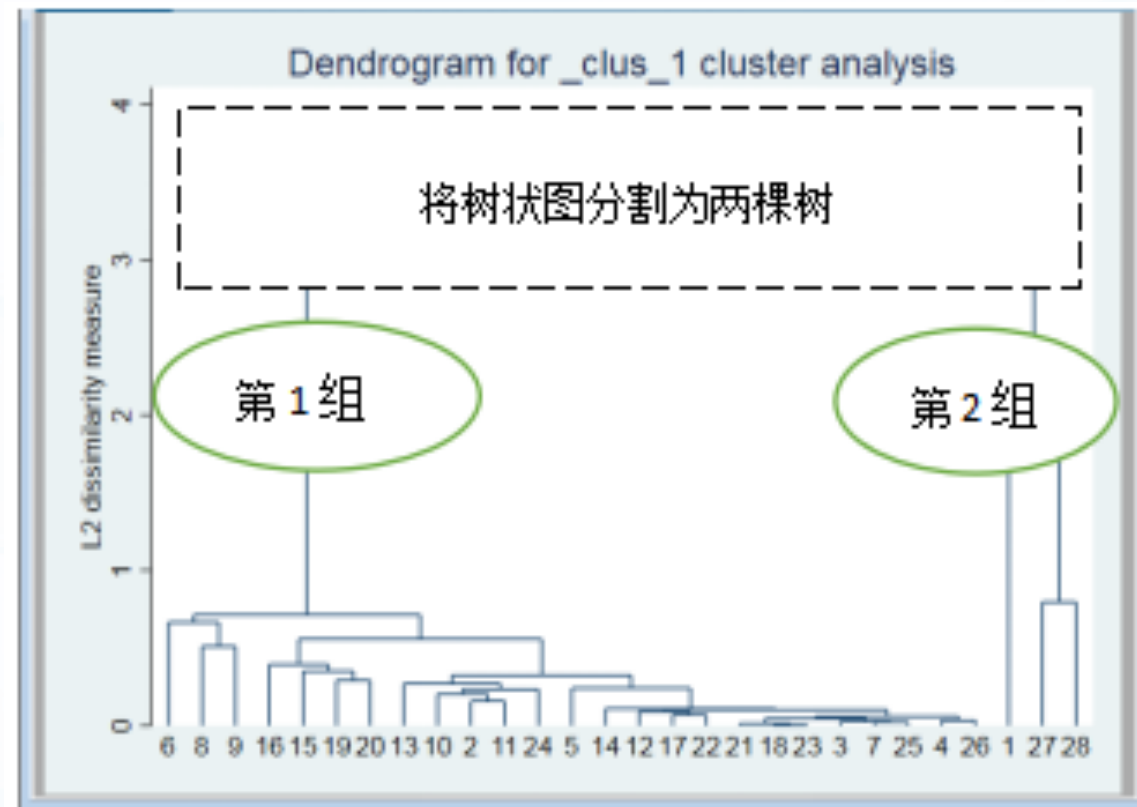
在实际分析中常用到的一个工具是树状图。图12.1即为树状图的一个示例。（图片来源：《Stata统计分析从入门到精通》杨维忠、张甜著，清华大学出版社，第12章聚类分析）。

观察该图，可以直观地看到具体的聚类情况，如21、18、23号样本聚合在一起（当然图片放大后，可能会进一步细分，比如21、18、23号样本首先是18、23号合并，然后与21合并），3、7、25号样本聚合在一起。



|| 层次聚类分析

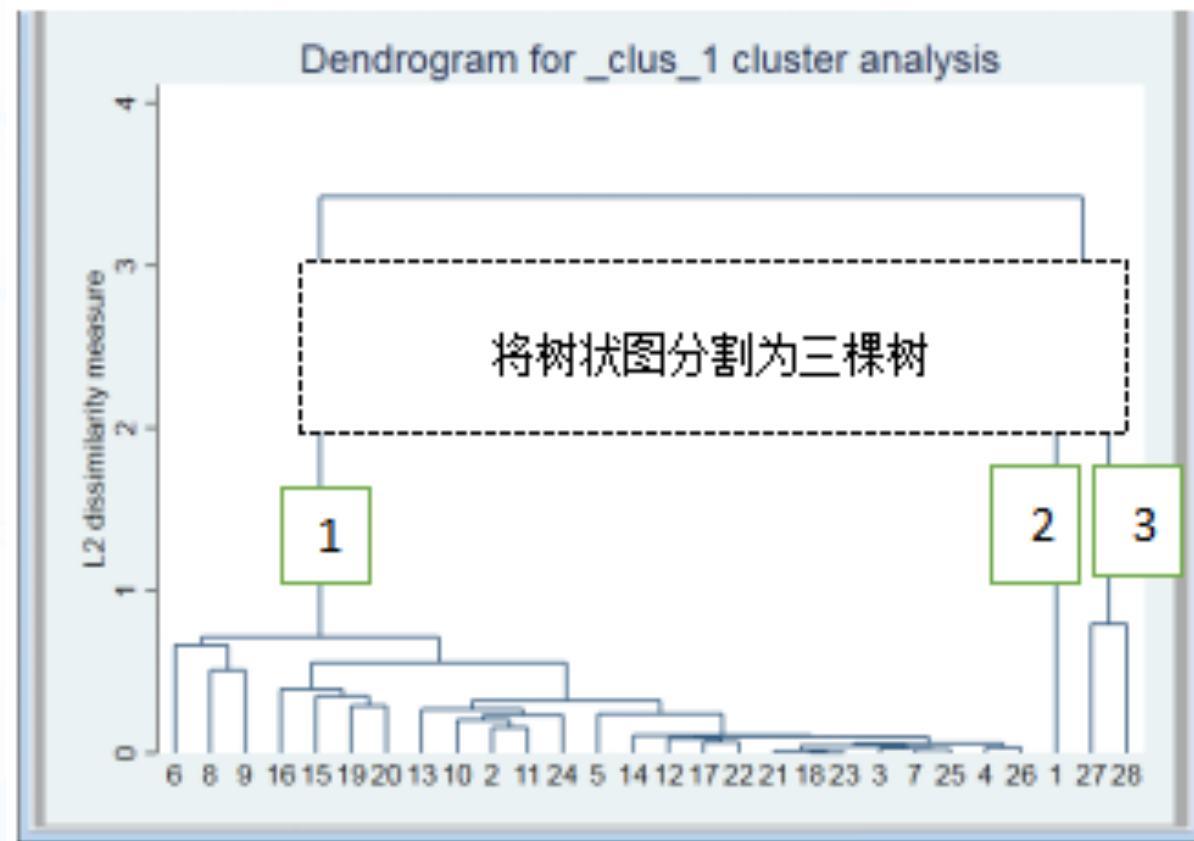
那么，到底分成了多少类呢？取决于研究的需要和实际的情况，需要用户加入自己的判断，确定好分类需求后，从聚类分析树状图最上面使用分割框往下进行分割。例如需要把所有样本观测值分成2类时，分析结果如图12.2所示。



|| 层次聚类分析

可以发现1、27、28三个样本为第2类，其他样本为第1类。如果需要把所有样本观测值分成3类，就需要自上而下继续进行分割（图中的分割框下移），分析结果如图12.3所示。

可以发现1号样本为第2类，27号、28号样本为第3类，其他样本为第1类。



层次聚类分析

与划分聚类分析方法相比，层次聚类分析方法的计算过程更为复杂，计算速度相对较慢，但是它不要求事先指定需要分类的数量，这一点是符合聚类分析探索性的本质特点的，所以这种聚类分析方法应用也非常广泛。

氏■■■技螯■颀■	■蔽■箔■■■■■掊■颀哏
跬■■■■■技螯 ° Single-Linkage Cluster Analysis `	帜搏■蔽■箔■■■■■柄变■表跬■■■貌■涤箔■■■■■ Y
跬■■■■■技螯 ° Complete-Linkage Cluster Analysis `	帜搏■蔽■箔■■■■■柄变■表跬■■■貌■涤箔■■■■■ Y
薏寿■■■■■技螯 ° Average-Linkage Cluster Analysis `	帜搏■蔽■箔■■■■■柄变■表龄■貌■涤■薏寿竦箔■■■■■ Y
表大竦■■■■■技螯 ° Median-Linkage Cluster Analysis `	帜搏■蔽■箔■■■■■柄变■表龄■貌■涤■表大竦箔■■■■■ Y
■徽■■■■■技螯 ° Centroid-Linkage Cluster Analysis `	帜搏■蔽■箔■■■■■柄变■表龄■貌■涤■■■徽箔■■■■■ Y
ward ■■■■■技螯 ° ward-Linkage Cluster Analysis `	变标■■■箔■■■■■拆■兼■徽■粉莞薏颀冒！登■瑚龄■貌■涤■粉莞薏颀冒跬综！■■■粉莞薏颀冒匕料■香 Y

$$\text{dist}_{xy} = \sqrt{\sum_{u=1}^n (x_u - y_u)^2}$$

衡量样本示例距离的测度

无论是划分聚类分析还是层次聚类分析，聚类分析的本质都是按照一定的距离对样本示例进行分割，在一定标准距离范围以内的样本示例被划分为一个类别，所以如何测度样本数据之间的距离也非常重要。

常用的针对连续变量数据的标准如表所示。

简写	含 义	适用变量类型
Euclidean ·L2 ·	<div>□□□□Y □□恨哏柄龄□啡箔□蕙颀堯箔冒□蕙颀□</div> $\text{dist}(X,Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$	□□粮□竦腌
L2squared	<div>□□□□□蕙颀</div>	□□粮□竦腌
Manhattan ·L1 ·	<div>跽渺□□□ !擦□□着啡□□</div>	□□粮□竦腌
Linfinity	<div>跽香啡□□</div>	□□粮□竦腌
correlation	<div>□哄□竦□粹忸掊□</div>	□□粮□竦腌
L(#)	<div>蛄料乡聃绦□□ !□□尿管磷瑚焚帜 p 啡 Y □□恨哏柄龄□啡箔□ p □蕙□着堯箔冒□ p □□</div> $d_{xy} = \left(\sum_{u=1}^n x_u - y_u ^p \right)^{\frac{1}{p}}$ <div>p 啡炼啡□射柄 1~7Y 喃 p=1 蜉 !丐蟪 书 [° 跽渺□□□ ≅喃 p=2 蜉 !丐蟪□哏□□ ≅喃 p □□拆蚯□香蜉 !丐蟪冀□□乡□□</div>	□□粮□竦腌
cosine	<div>怠喋□粹掊</div>	□□粮□竦腌

|| 衡量样本示例距离的测度

常用的针对分类变量数据的标准如表所示。

针对聚类分析算法，同样要求首先对特征变量数据进行标准化，之所以这么做，是因为如果不进行标准化，少数变量可能会对距离的影响太大，使得分析结果严重失真。

简写	含 义	适用变量类型
matching	<input type="checkbox"/> 狂空 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Jaccard	Jaccard <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Russell	Russell 冒 Rao <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Hamann	Hamann <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Dice	Dice <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
antiDice	镰 Dice <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Sneath	Sneath 冒 Sokal <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Rogers	Rogers 冒 Tanimoto <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Ochiai	Ochiai <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Yule	Yule <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Anderberg	Anderberg <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Kulczynski	Kulczynski <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Pearson	Pearson <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌
Gower2	蔽 Pearson <input type="checkbox"/> 领技 <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> 粹 <input type="checkbox"/> 竦	技 <input type="checkbox"/> 粮 <input type="checkbox"/> 竦腌



PART 02

数据准备

数据准备

本节我们以“数据12.1”为例进行讲解，其中的数据为《中国2019年分地区连锁餐饮企业基本情况统计》，摘编自《中国统计年鉴2020》。这个数据文件中共有9个变量，分别是V1~V9，分别表示地区、总店数、门店总数、年末从业人数、年末餐饮营业面积、餐位数、营业额、商品购进总额、统一配送商品购进额。

下面我们以V3门店总数、V6餐位数、V7营业额、V8商品购进总额4个变量对所有样本观测值开展划分聚类分析。

载入分析所需要的模块和函数

在进行分析之前，我们首先载入分析所需要的模块和函数，读取数据集并进行观察。

示例

[参阅教材内容](#)

|| 变量设置及数据处理

示例

参阅教材内容

|| 特征变量相关性分析

示例

参阅教材内容



PART 03

划分聚类分析算法示例

|| 使用K均值聚类分析方法对样本示例进行聚类(K=2)

示例

参阅教材内容

|| 使用K均值聚类分析方法对样本示例进行聚类(K=3)

示例

参阅教材内容

|| 使用K均值聚类分析方法对样本示例进行聚类(K=4)

示例

参阅教材内容



PART 04

层次聚类分析算法示例

最短联结法聚类分析

示例

参阅教材内容

|| 最长联结法聚类分析

示例

参阅教材内容

|| 平均联结法聚类分析

示例

参阅教材内容

|| ward联结法聚类分析

示例

参阅教材内容

|| 重心联结法聚类分析

示例

参阅教材内容



PART 05

习 题

习题

继续使用“数据12.1”，以总店数、年末从业人数、年末餐饮营业面积、统一配送商品购进额4个变量，即V2、V4、V5、V9，对所有样本观测值开展划分聚类分析和层次聚类分析。

- (1) 载入分析所需要的库和模块
- (2) 变量设置及数据处理
- (3) 特征变量相关性分析
- (4) 使用K均值聚类分析方法对样本示例进行聚类(K=2)
- (5) 使用K均值聚类分析方法对样本示例进行聚类(K=3)

(6) 使用K均值聚类分析方法对样本示例进行聚类($K=4$)



感谢聆听

THANKS
