

# 第十四章 数据探查建模实验

通过数据处理，挖掘其背后隐藏的信息是数据处理的根本目的。根据数据处理的任務，大体分类两类，一类是以探查数据的特征，寻找数据符合某种规律为任务的数据处理称为数据探查建模；一种是根据已知样本数据特征，构建数据模型，实现对未知样本数据的预测称为数据预测建模。也就是获得一批数据后，希望对这批数据寻找到某种规律性，就属于探查建模，已知一批数据的规律性，希望对于未知数据进行预测，这就是预测建模。它们也称为无模型训练和有模型训练。本章实验主要讨论探查建模实验，包括聚类分析和关联分析等，编写过程中参考了文献[20, 22]和文献[23]，也参考了网上的一些短文材料。

## 第一节 聚类分析

获得一批数据，对这批数据进行简化描述的最好手段就是聚类，也就是根据数据的指标特征，分析哪些数据具有相似的特征，从而归为一类，这样就将纷繁复杂的数据通过简单的归类进行简化，同时也把握了其本质特征。比如刚入学的大学新生，学校希望进行归类，比如根据英语成绩、数学成绩和兴趣爱好等进行归类，也就是将有相同特长和兴趣爱好的学生进行归类，然后开展有针对性的教学，这就是聚类分析。

根据聚类是否指定聚类数目，可以将聚类分为两种，一种是层次聚类，一种是划分聚类。层次聚类就是最初每一个样本都是一类，然后根据给定标准进行两两合并，一直到所有的样本都聚为一类为止；而划分聚类是根据给定聚类标准，将样本归类为给定数目。例如有各种鱼、禽和肉类食品，如果根据食品中的一些营养标准，就可以这些食品聚为几个类别。

### 一、聚类的步骤

对于一批数据进行聚类，显然需要讨论聚类的标准，没有标准就没有办法进行归类，因此在聚类的时候需要对数据进行一系列的处理。

#### 1. 选择合适的变量

我们对于一批样本，它可能有许多指标取值，比如上面谈到的鱼和肉等食品，它们所含的营养指标：氨基酸，维生素，蛋白质，糖类，铁和钠等含量。每一个样本都在这些指标上取值，如果我们对这些样本进行聚类，目的是关注营养成分，可能我们只能选择氨基酸和蛋白质以及糖类等作为聚类指标，舍弃其它非营养成分的含量。

#### 2. 数据标准化

由于不同指标，量纲不同，有的指标可能取值范围为 0 到 1 之间，有的取数百万大小，这样一起聚类，可能取值较大的指标淹没了取值较小的指标，造成聚类效果不佳，因此可以统一使用标准化等手段，比如 R 语言的 `scale()` 函数，将所有指标取值都标准化到 0 和 1 之间；

#### 3. 异常点的处理

某些指标可能由于误差或者其它因素，取值异常，这就可能对聚类效果造成干扰。R 语言包 `outliers` 中的函数 `mvoutlier()` 可以帮助识别离群点；

#### 4. 计算样本之间的距离

两个样本是否可以归为一类，总要有一个标准，这个标准就是衡量两个样本之间的距离大小，距离越大越不能归为一类，因此这个距离的选取十分关键和重要，目前有很多距离计算方法，比如欧氏距离，非对称二元距离等等，R 语言有一个 `dist()` 函数，提供很多计算指标距离的方法

#### 5. 聚类算法的确定

是选择层次聚类还是选择划分聚类，不同的方法有不同的适应范围，比如层次聚类适合小样本的聚类。

#### 6. 确定类的数目

最终聚类成几类，需要确定类的数目，如果给定一个标准，也就是指标距离足够大，那么所有的样本都可能归为一类。因此给定距离或者给定聚类数目，都可以将给定样本进行归类。

## 7. 可视化结果

将聚类结果做可视化输出，比如层次聚类就可以使用树状图显示出来

## 8. 聚类结果的解读

对于聚类结果，可以观察哪些样本聚为一类，比如根据一些营养指标，将某些鱼类和禽类聚为一类，那么我们可以解读出，这几种食品具有相同的营养价值。

## 9. 验证结果

验证结果讨论的是聚类的稳定性问题，也就是给一批同样性质的样本，它们是否也会输出相近似的聚类结果。如果不能，就说明聚类模型不可靠，需要改进。R 语言的 `fpc`、`clv` 和 `clValid` 等包提供了一些评估聚类解函数稳定性的验证方法。

## 二、层次聚类分析

这种聚类方法是通过迭代步骤实现的。刚开始所有的样本归为一类，然后计算两两样本之间的距离，距离最小的聚为一类，这样得到新的样本分类，然后再每两个类计算距离，将距离小的再聚为一类，这样持续下去直到所有的样本都聚为一类为止。这里面要注意的是如果两个样本已经聚为一类了，那么就必须立刻作为一个整体参与和其它样本或者样本类计算距离。

这里面就有一个问题，那就是计算两个样本之间距离很简单，比如使用欧氏距离就可以了，但是计算两个类之间的距离怎么办呢？很显然也要通过计算一个类的样本和另外一个类的样本之间的距离来表示。这时候就有很多方法，比如单联动计算方法就是两个类中样本的最小距离；全联动法就是一个类中样本和另外一个类中样本的最大距离；平均联动就是两个类所有点之间的平均距离；质心法就是两个类质心的距离；Ward 法：就是两个类之间所有变量的方差的平方和。

层次聚类方法在 R 语言中可以使用 `hclust` 函数来实现，这个函数的格式是 `hclust(d,method=)`，其中 `d` 是通过 `dist()` 函数生成的距离矩阵，`method` 可以取值为 `single`、`complete`、`average`、`centroid` 和 `ward` 等。对应的就是上面的单联动方法和全联动方法等等。R 语言中有一个包 `flexclust`，该包有一个数据集为 `nutrient`，是化验各种食品的热量、蛋白质、脂肪、含钙量和含铁量得到的数据样本，现在要对这些样本进行聚类，具体 R 实现代码如下，柱状聚类图如图 14.1 所示。

```
install.packages("flexclust")
data(nutrient, package="flexclust") # 这条语句就是装载包 flexclust 包的数据集
row.names(nutrient) = tolower(row.names(nutrient)) # 将 nutrient 中的字体变成小写字体
nutrient.scaled = scale(nutrient) # 对数据集进行标准化
d = dist(nutrient.scaled) # 默认采用欧氏距离
fit.average = hclust(d, method="average") # 使用平均联动方法进行聚类
plot(fit.average, hang=-1, cex=.8, main="Average Linkage Clustering")
```

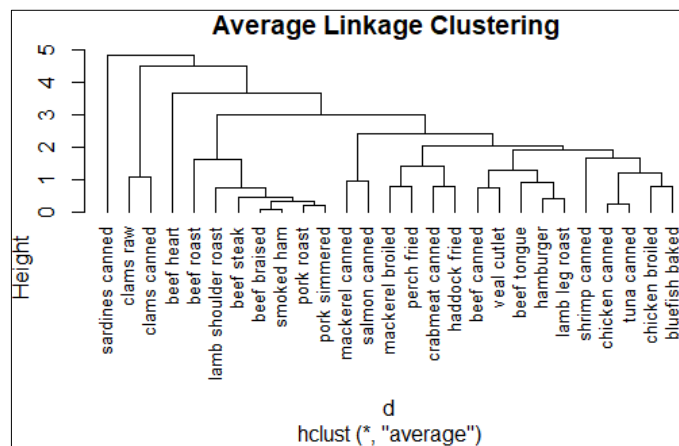


图 14.1 食品数据的平均联动聚类

图 14.1 从下面向上看，我们看所有食品中，距离最小的两个样本是 `beef braised` 和 `smoked`

ham，它们当然合并在一起，接着计算发现 pork roast 和 pork simmered 距离最小，因此合并为一类，这两类和其它所有样本再两两计算距离，最小的在合并为一类，依次进行下去，就会得到如图 14.1 所示的树状图。注意计算两个类之间的距离使用的平均联动方法。分析图 14.1，我们可以观察哪些食品营养价值是较为相似的，哪些存在很大不同，进而给客户购买意见。

### 三、划分聚类分析

划分聚类分析主要分为两种，一种是 K-均值聚类，一种是基于中心点的划分聚类。

#### 1. K 均值聚类

##### (1) K 均值聚类具体步骤

随机选择 K 个样本点作为中心点；将每一个样本点归类到距离它最近的中心点；重新计算得到的 K 个类的中心点；将每一个样本点重新归类到距离它最近的中心点；重复上面的步骤，直到样本不再被重新分配或者达到最大迭代次数。

##### (2) K 均值聚类的优缺点

K 均值可以处理比层次聚类更大的数据集，缺点是要求变量是连续变化的，同时受异常值影响大，不太适用于非凸数据集。

##### (3) R 语言的实现

由于 K 均值聚类需要指定聚类数目，指定不合适的聚类数目会使得聚类效果很差，因此 R 语言的包 NbClust 提供了 NbClust() 用来探查数据究竟分为几类较为合适。NbClust() 的调用形式为 NbClust(data, distance=, min.nc=, max.nc=, method=)，其中 distance 可以选择 "euclidean"，method 可以选择 "average"。NbClust 提供用来衡量数据集聚为 n 类的指标数，下面例子可以看到类别和指标支持数目。

下面以 R 语言包 rattle 中的 wine 数据集来实验 K 均值聚类，下面是具体 R 语言代码

```
install.packages("NbClust")
install.packages("rattle")
data(wine, package="rattle") # 将数据集装载到内存
df = scale(wine[-1])        # 去除第一列数据，然后标准化
library(NbClust)           # 装载包
set.seed(1234)             # 再次运行该代码时候，随机数和上次相同，如果没有这个语句，
                           # 聚类的类号编号可能不是按照 123 的顺序编号，这个是因为样本类别编号产生
                           # 的随机数也是按照此设置的，这样保证聚类号和样本号一致
nc = NbClust(df, min.nc=2, max.nc=15, method="kmeans")
barplot(table(nc$Best.nc[1,]), xlab="Number of Clusters", ylab="Number of Criteria")
```

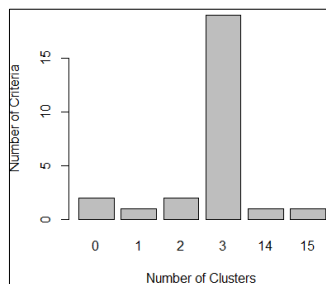


图 14.2 酒类数据集推荐的聚类个数

具体代码含义上面都已经注释了，图 14.2 是生成的聚类数目和支持指标数目的柱状图，通过该图观察，NbClust 这个包针对 wine 数据集提供了 26 种指标，在该 26 种评判指标中，有超过 14 种指标建议选择聚为 3 类最合适。代码中使用 table(nc\$Best.nc[1,]) 获得支持聚类数目的指标数。可以使用 t(nc\$Best.nc) 查看 26 种具体指标名称。下面是使用 R 语言基础包提供的 K 均值聚类函数，选择聚为 3 类针对 wine 数据集进行聚类的 R 语言代码。

```
set.seed(1234)
fit = kmeans(df, 3, nstart=25)
fit$size
fit$centers
```

变量 fit\$size 存放的是 65, 51 和 62，分别表示三个类的包含样本个数。变量 fit\$centers 存放的是三个类的中心点，显然一个中心点就是该类样本的中心。由于 wine 数据集包含 13 个指标，

因此很难可视化的展示出来，如果只包含两个指标，则可以使用二维坐标平面可视化出来，三个指标可以使用三维坐标可视化出来。变量 `fit$cluster` 存放的是每一个样本被归为三个类中的哪一个类，或者说哪些样本聚为一类了。另外注意的一点就是由于聚类的数据都标准化了，因此计算的中心点也是标准化后的值，可以使用 `aggregate()` 函数将标准化的中心点坐标恢复到原始数据的形式。具体 R 语言代码如下

```
aggregate(wine[-1], by=list(cluster=fit$cluster), mean)
```

数据集 `wine` 的第一列是样本实际上所属的类别，可以将它和 K 均值聚类比较，观察 K 均值聚类的效果，比如使用 `ct = table(wine$Type, fit$cluster)` 生成列联表如下

	1	2	3
1	59	0	0
2	3	65	3
3	0	0	48

可以观察到有 6 个样本分类错误，也就是其中第 2 类有 3 个样本分错到 1 类中，有 3 个第 2 类的样本分错到 3 类中。可以使用 `str()` 函数查看 `wine` 数据集中每一个样本含有的 13 个成分是什么物质，比如酒精含量，钙的含量，黄酮素的含量等等。聚为三类意味着这些酒根据成分可以划分为三类。

## 2. 基于中心点的划分聚类(PAM)

K 均值聚类是基于样本均值聚类，也就是根据距离样本中心点的距离来聚类，很显然如果一个非常大的异常点，就会使得中心点严重偏离大多数样本点，这就使得 K 均值聚类变得不稳定了，因此提出基于中心点的划分聚类。

PAM 与 K-means 相比，PAM 可以使用任意度量来计算两个样本的距离，因此可以容纳混合数据，不限于连续变量，另外就是 PAM 算法和 K-means 算法根本区别是 PAM 选择一个最具有代表性的样本来代替 K-means 的中心点，因此 PAM 稳健性更强。很显然选择代表性的点需要不停的迭代计算，具体算法的思想就是不再使用 K-means 的中心点计算，而是任意选一个点作为中心点，然后计算所有点和该点距离的和，再选另外一个点作为中心点，再计算这种距离和，比较大小，距离和小的点作为代表点，重复计算下去，直到收敛或者达到指定次数为止。通过这种算法，可以观察到 PAM 计算量很大，收敛速度较慢。

使用 PAM 对上一节的 `wine` 数据集进行聚类可以使用 `cluster` 包中的 `pam()` 函数，注意 `kmeans()` 是 R 语言基础包的函数，而使用 `pam()` 函数需要加载 `cluster` 包，具体 R 语言代码如下

```
install.packages("cluster")
data(wine, package="rattle")
library(cluster)
set.seed(1234)
fit = pam(wine[-1], k=3, stand=TRUE) # k 表示分为 3 类，stand 是一个逻辑值，表示聚类时是否需要标准化
```

我们可以观察返回结果 `fit` 中包含的中心点，样本所属的类等等。同样也可以使用 `table()` 函数生成列联表即：`table(wine$Type, fit$clustering)`，生成的列联表如下：

	1	2	3
1	59	0	0
2	16	53	2
3	0	1	47

可以观察到，第二类中有 16 个样本被错分到第 1 类中了，有 2 个被错分到第 3 类中了。第 3 类有一个样本错分到第 2 类中。这表明 PAM 聚类对于这类数据集不如 k-means 聚类更有效。

## 3. 划分聚类的可视化

划分聚类对于变量数目超过 3 个的数据集，很难可视化表示出聚类的情况，但是如果使用主成分分析方法，将最大的两个主成分指标作为平面坐标系的两个坐标轴，倒是可以可视化聚类的效果，R 语言提供了这么一个函数 `clusplot()` 完成该功能，具体实现很简单就是 `clusplot(fit)`，其中 `fit` 是聚类的结果，使用该函数绘制的 PAM 对 `wine` 数据集聚类的结果如图 14.3 所示，注意其中的椭圆是人为绘制的，是包含同一类样本的最小椭圆。对于 K-means 聚类的结果，可以使用 `clusplot(df, fit$cluster)` 查看可视化聚类结果。注意 `clusplot()` 函数是 `cluster` 包中的函数，如果没有转载入该包，需要运行 `library(cluster)`，否则运行会出错。

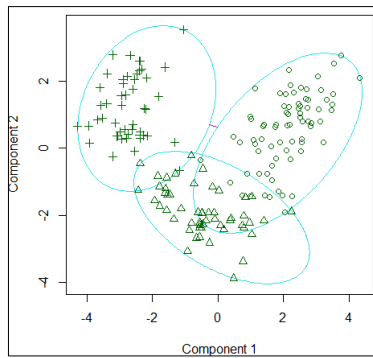


图 14.3 酒类数据集 PAM 聚类的两个主成分可视化图形

## 第二节 关联分析

关联分析就是计算事务中一些项同时出现的频率。这在现实生活中很常见，比如图书馆借书，每个学生一次借书就是一个事务，一次借书事务中一共借了几本书，就是几项。关联分析就是在一段时间内，分析哪些书是一起借阅的，了解这些信息就可以为图书馆里人员排列书架等提供依据，因此关联分析也是一种从数据集中探查建模的数据处理方法，具体就是探查出事务数据集中的关联规则，寻找哪些项目是有联系的。关联分析寻找事务集中的强关联规则通常使用 Apriori 算法，其中需要两个指标，一个是支持度，一个是置信度，比如对于关联规则  $X \rightarrow Y$ ，支持度表示同时包含  $X$  和  $Y$  的事务占总事务的比例，置信度指的是同时包含  $X$  和  $Y$  的事务在所有包含  $X$  事务的占比。通过支持度和置信度来确定关联规则。

R 语言有数据包 `arules` 的 `apriori()` 实现关联分析。该数据包中还有 `Groceries` 数据集可供关联分析做实验，具体 R 语言代码如下：

```
install.packages("arules")
library(arules)
data(Groceries)
groceryrules = apriori(Groceries, parameter = list(support = 0.006, confidence = 0.25, minlen = 2))
inspect(groceryrules[1:5])
```

其中 `inspect()` 函数用来观察得到的关联规则，比如查看部分关联规则如下

	lhs	rhs	support	confidence	coverage	lift	count
[1]	{pot plants}	=> {whole milk}	0.006914082	0.4000000	0.01728521	1.565460	68
[2]	{pasta}	=> {whole milk}	0.006100661	0.4054054	0.01504830	1.586614	60
[3]	{herbs}	=> {root vegetables}	0.007015760	0.4312500	0.01626843	3.956477	69
[4]	{herbs}	=> {other vegetables}	0.007727504	0.4750000	0.01626843	2.454874	76
[5]	{herbs}	=> {whole milk}	0.007727504	0.4750000	0.01626843	1.858983	76

可以发现这些规则是根据最小支持度 0.006，最小置信度 0.25 生成的，`inspect(groceryrules[1:5])` 返回前 5 个关联规则，可以发现第一个关联规则就是盆栽植物和全脂牛奶有一定的关联度，这意味着有一批客户在买盆栽植物时候也一并买了全脂牛奶，其支持度为 0.0069，置信度为 0.017。需要注意返回的关联规则并没有按照支持度或者置信度从小到大的顺序给出。

关联分析只能针对类别型的数据探查关联规则，连续型数据需要离散化，函数 `apriori()` 中的第一个参数必须是事务型数据结构才能使用该函数。下面举例一般事务集如何进行关联分析。

1. 生成 `list` 数据结构，也就是将所有事务放入一个 `list` 数据结构。每一个事务就是一个向量，向量的每一个元素就是一项，是使用双引号引起来的字符串。例如下面生成的事务集 `alist`。

```
alist = list(
  c("a", "b", "c"),
  c("a", "b"),
  c("a", "b", "d"),
  c("c", "e"),
  c("a", "b", "d", "e"))
```

2. 给每一个事务指定一个名称。比如下面语句使用 `paste()` 函数生成一个字符串向量，其中 `Tr` 和后面的向量分别组合一个字符串，一共 5 个字符串，这 5 个字符串组成一个字符串向量。

```
names(alist) = paste("Tr",c(1:5), sep = "") # 作用就是给每一个事务一个名称，方便生成关联规则
```

### 3. 生成一个事务型结构的变量

```
trans1 = transactions(alist)
```

### 4. 进行关联分析和显示

```
rules = apriori(trans1) #使用默认的支持度和置信度  
inspect(rules)
```

通过以上几个步骤，即可以获得关联规则，对于以数据文件存放的数据，同样可以使用 `readTable()` 函数读取获得的 `dataframe` 数据结构数据，然后使用 `transactions()` 函数转换为事务型数据变量

## 第三节 综合实验十四

### 一、实验目的

1. 了解数据探查建模的基本思想
2. 掌握聚类分析中 K-means 和 PAM 聚类的 R 语言实现，理解相关函数中的各种参数
3. 掌握关联分析的 R 语言实现，能够具体应用

### 二、实验平台和软件

1. 互联网平台
2. 计算机系统
3. Rstudio 环境和 R 语言编译系统

### 三、实验内容和步骤

#### 1. 层次聚类分析

- (1) 获取 `flexclust` 包中的 `nutrient` 数据集
- (2) 阅读 `nutrient` 数据集的特点
- (3) 对该数据集进行标准化并计算样本距离
- (4) 对计算的距离矩阵进行层次聚类
- (5) 绘制层次聚类图
- (6) 分析该聚类结果的含义

#### 2. K 均值聚类

- (1) 获取包 `rattle` 中的 `wine` 数据集
- (2) 了解该数据集的特点
- (3) 标准化该数据集
- (4) 使用 `NbClust()` 函数探查该数据集聚几类最为合适
- (5) 使用 `Kmeans()` 函数进行聚类
- (6) 将聚类结果和实际分类使用 `table()` 函数生成列联表
- (7) 了解 K-means 聚类的效果

#### 3. PAM 聚类

- (1) 安装和装载包 `cluster`
- (2) 使用 `pam()` 函数对标准化的 `wine` 进行聚类
- (4) 将聚类结果和实际分类使用 `table()` 函数生成列联表
- (5) 了解 PAM 聚类效果

#### 4. 关联分析

- (1) 安装和装载包 `arules` 包
- (2) 获取 `Groceries` 数据集
- (4) 使用 `apriori()` 函数对该数据集进行关联分析

(5) 使用 inspect() 函数了解关联规则

(6) 分析和解释一些规则的含义

#### 四、实验报告 (总分 8 分)

1. R 语言中有 iris 数据集，分析说明该数据集各指标含义，探查聚类的合理数目，然后使用 K-means 进行聚类，并分析聚类的有效性，最后可视化聚类结果。(2 分)

2. 使用层次聚类方法对 iris 进行聚类，绘制层次聚类图，并和 K-means 聚类的结果进行比较，分析各自特点。(2 分)

3. 下面是某个超市购物记录，使用关联分析探查一些关联规则，并给出支持度和置信度。(4 分)

交易时间	购买商品
14: 25	i1, i2, i4
15: 07	i1, i2, i3
16: 33	i2, i3
17: 05	i1, i3
18: 40	i1, i2, i3, i5
18: 55	i2, i3
19: 26	i1, i2, i5
19: 52	i2, i4
20: 03	i1, i2, i3
20: 16	i1, i3