

周志华 著

MACHINE  
LEARNING

机器学习

清华大学出版社

## 第9章 聚 类

### 9.1 聚类任务

常见的无监督学习任务还有密度估计 (density estimation)、异常检测 (anomaly detection) 等。

对聚类算法而言, 样本簇亦称“类”。

聚类任务中也可使用有标记训练样本, 如 9.4.2 与 13.6 节, 但样本的类标记与聚类产生的簇有所不同。

在“无监督学习”(unsupervised learning)中, 训练样本的标记信息是未知的, 目标是通过对无标记训练样本的学习来揭示数据的内在性质及规律, 为进一步的数据分析提供基础。此类学习任务中研究最多、应用最广的是“聚类”(clustering)。

聚类试图将数据集中的样本划分为若干个通常是不相交的子集, 每个子集称为一个“簇”(cluster)。通过这样的划分, 每个簇可能对应于一些潜在的概念(类别), 如“浅色瓜”“深色瓜”, “有籽瓜”“无籽瓜”, 甚至“本地瓜”“外地瓜”等; 需说明的是, 这些概念对聚类算法而言事先是未知的, 聚类过程仅能自动形成簇结构, 簇所对应的概念语义需由使用者来把握和命名。

形式化地说, 假定样本集  $D = \{x_1, x_2, \dots, x_m\}$  包含  $m$  个无标记样本, 每个样本  $x_i = (x_{i1}, x_{i2}, \dots, x_{in})$  是一个  $n$  维特征向量, 则聚类算法将样本集  $D$  划分为  $k$  个不相交的簇  $\{C_l \mid l = 1, 2, \dots, k\}$ , 其中  $C_l \cap C_{l'} = \emptyset$  且  $D = \bigcup_{l=1}^k C_l$ 。相应地, 我们用  $\lambda_j \in \{1, 2, \dots, k\}$  表示样本  $x_j$  的“簇标记”(cluster label), 即  $x_j \in C_{\lambda_j}$ 。于是, 聚类的结果可用包含  $m$  个元素的簇标记向量  $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_m)$  表示。

聚类既能作为一个单独过程, 用于找寻数据内在的分布结构, 也可作为分类等其他学习任务的前驱过程。例如, 在一些商业应用中需对新用户的类型进行判别, 但定义“用户类型”对商家来说却可能不太容易, 此时往往可先对用户数据进行聚类, 根据聚类结果将每个簇定义为一个类, 然后再基于这些类训练分类模型, 用于判别新用户的类型。

基于不同的学习策略, 人们设计出多种类型的聚类算法。本章后半部分将对不同类型的代表性算法进行介绍, 但在此之前, 我们先讨论聚类算法涉及的两个基本问题——性能度量和距离计算。

### 9.2 性能度量

聚类性能度量亦称聚类“有效性指标”(validity index)。与监督学习中的