

聚类分析算法简介

探索无监督学习中的核心技术,将数据按相似性智能分组



什么是聚类分析算法？

聚类分析算法是一类**无监督学习**方法,它的目标是将数据集中的对象按照相似性分组。这种算法不需要预先标注的训练数据,而是通过发现数据内在的结构模式来自动分类。

聚类的核心理念是让相似的对象聚在一起,不相似的对象分开,从而揭示数据中隐藏的群组结构。

组内相似性最高

同一类中的对象尽可能相似

组间差异性最大

不同类之间的对象尽可能不同

聚类算法的工作原理



()	([
数据预处理 标准化、归一化处理	相似性度量 定义什么叫"相似"
({	(}
聚类策略 采用何种方法分组	结果验证 评估分类效果

经典聚类算法



K-means算法:最流行的聚类方法



基本思想

将数据分成 k 个簇,每个簇用其中心点代表。算法通过迭代优化,不断调整中心点位置,直到达到最优分组。

生活类比

像在地图上选 k 个商业中心,每个人去最近的商业中心购物,然后根据客流重新调整商业中心位置。

K-means算法步骤



选择初始中心点

随机选择k个数据点作为初始簇中心



重新计算中心

计算每个簇内所有点的平均值作为新中心



分配数据点

将每个数据点分配给距离最近的中心点



迭代优化

重复步骤2-3,直到中心点不再变化

K-means算法的优缺点

优点

- 简单易理解,计算效率高
- 适用于大规模数据集
- 结果相对稳定

缺点

- 需要预先指定k值
- 对初始中心点敏感
- 假设簇是球形的
- 对噪声和异常值敏感

适用场景



客户细分

市场分析与用户分组



图像处理

图像分割与数据压缩



推荐系统

用户兴趣分组

层次聚类算法:构建树状结构

层次聚类算法通过构建数据的树状层次结构来进行分组,不需要预先指定簇数。算法结果直观,能够展示数据在不同粒度下的分组关系。



自底向上

凝聚式:从每个点单独成簇开始,逐步合并相似的簇




自顶向下

分裂式:从所有点一个簇开始,逐步分割成更小的簇



层次聚类算法步骤(凝聚式)



 **生活类比:** 像画家族族谱,从每个人开始,根据关系远近逐步归并成家族分支。

层次聚类的特点与应用

优点

- 不需要预先指定簇数
- 能发现任意形状的簇
- 结果直观,有层次结构
- 结果确定,重复运行结果相同

缺点

- 时间复杂度高 $O(n^3)$
- 不适合大规模数据
- 对噪声敏感
- 难以处理高维数据

典型应用场景

物种分类

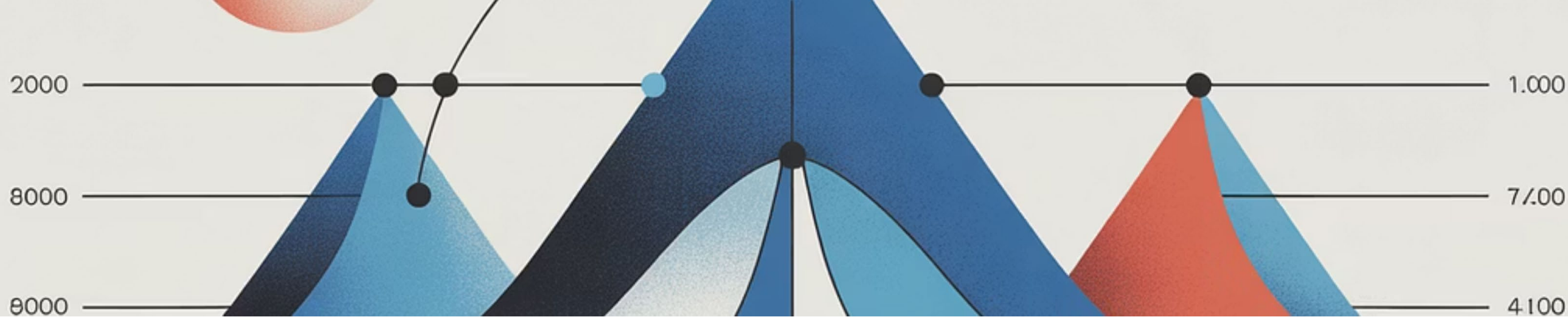
生物学中的物种分类与基因分析

社会网络

社交关系分析与社区发现

产品分类

构建产品分类体系



相似性度量方法

选择合适的相似性度量方法是聚类成功的关键。不同的距离度量适用于不同类型的数据和应用场景。

欧几里得距离

公式: $d(x,y) = \sqrt{[(x_1-y_1)^2 + (x_2-y_2)^2 + \dots + (x_n-y_n)^2]}$

最常用的距离度量,适用于连续数值变量

曼哈顿距离

公式: $d(x,y) = |x_1-y_1| + |x_2-y_2| + \dots + |x_n-y_n|$

对异常值不敏感,适用于高维稀疏数据

余弦相似度

公式: $\cos(x,y) = (x \cdot y) / (|x| \times |y|)$

衡量向量夹角,适用于文本分析和推荐系统

数据预处理的重要性

为什么需要标准化?

不同特征的量纲差异会严重影响聚类结果。例如:

- 年龄: 20-60
- 收入: 3000-50000

如果不标准化,收入的影响会远大于年龄,导致聚类结果偏向收入维度。

☐ 标准化将所有特征缩放到相同范围,确保每个特征对聚类的贡献相当。

缺失值处理策略



删除法

数据充足时直接删除含缺失值的记录



插值法

使用均值、中位数或K近邻插值填补

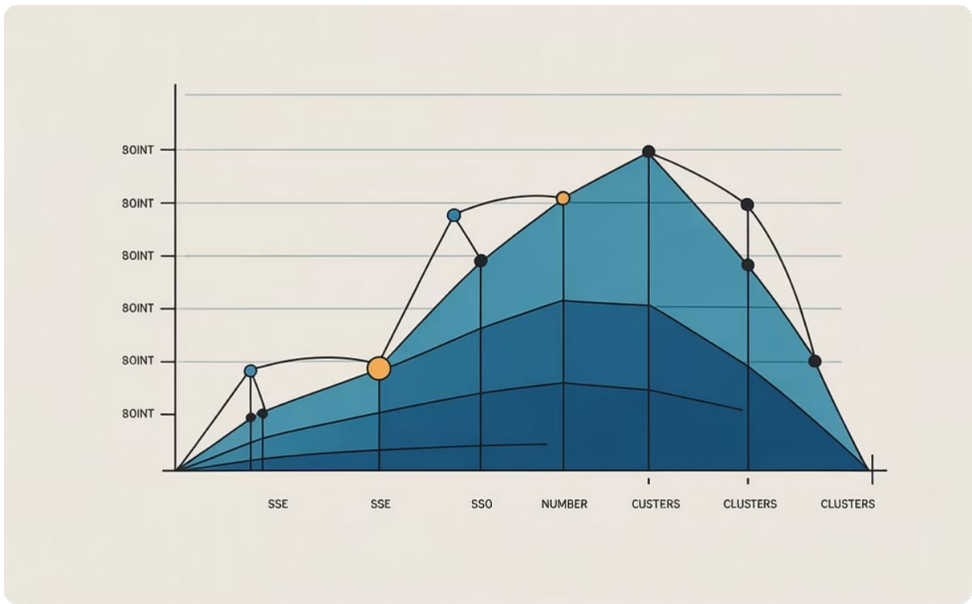


标记法

将缺失作为特殊类别处理

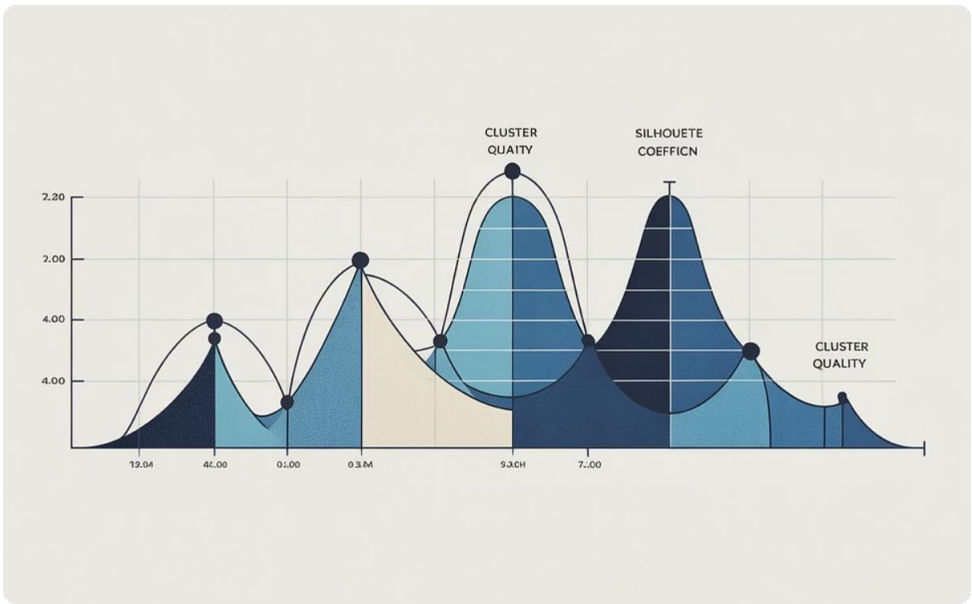
如何确定最佳簇数？

确定合适的簇数是聚类分析中的关键问题。以下是两种常用的方法：



肘部法则

绘制不同k值下的SSE(误差平方和)曲线,寻找曲线中的"肘部"拐点。拐点处代表增加簇数带来的收益开始递减,是较优的簇数选择。



轮廓系数

衡量簇内紧密度和簇间分离度的综合指标。取值范围为[-1,1],值越大表示聚类效果越好。通过比较不同k值的轮廓系数来选择最优簇数。

聚类分析的关键要点

算法选择

根据数据规模、形状特征和业务需求选择合适的聚类算法

数据预处理

标准化处理和缺失值处理是保证聚类质量的基础

相似性度量

选择适合数据特点的距离度量方法

参数优化

使用科学方法确定最佳簇数和其他关键参数

结果验证

通过多种评估指标验证聚类效果的合理性