



# Python

# 机器学习实战

# 聚类算法

## 聚类和分类的区别

聚类是无监督的学习

分类是有监督的学习

## 聚类的场景

使用聚类不需要提前被告知要划分的组是什么样的，  
在我们不知道找什么时就自动完成分组



## 无监督学习算法

### 聚类算法一览

层次聚类法

Kmeans

密度聚类法

谱聚类

GMM

LDA

... ..

## 基于距离

- 层次聚类法
- K-means聚类法
- K-center聚类法

## 基于密度

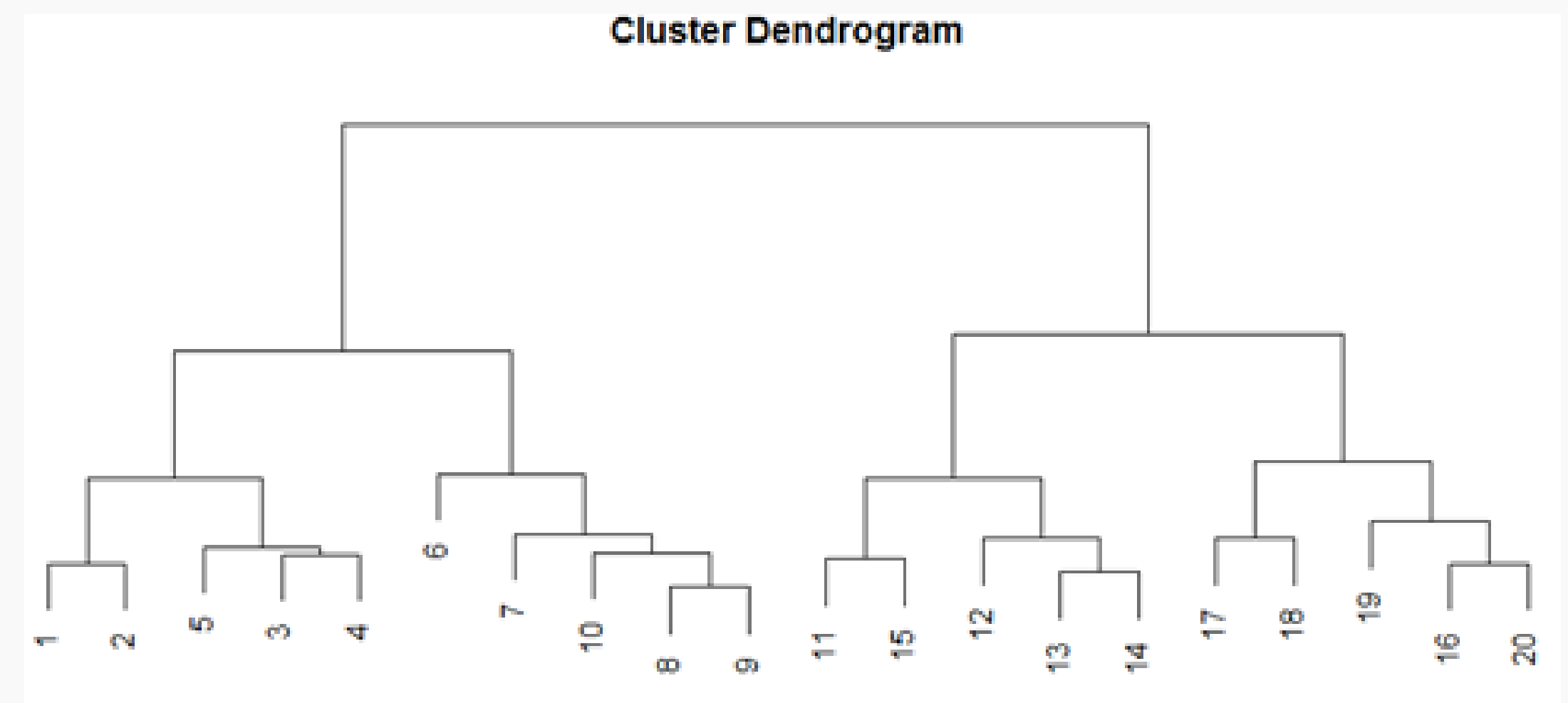
- 密度距离法 (DBCSAN)

## 数学中距离满足三个要求

- 必须为正数
- 必须对称
- 满足三角不等式

## 层次聚类法

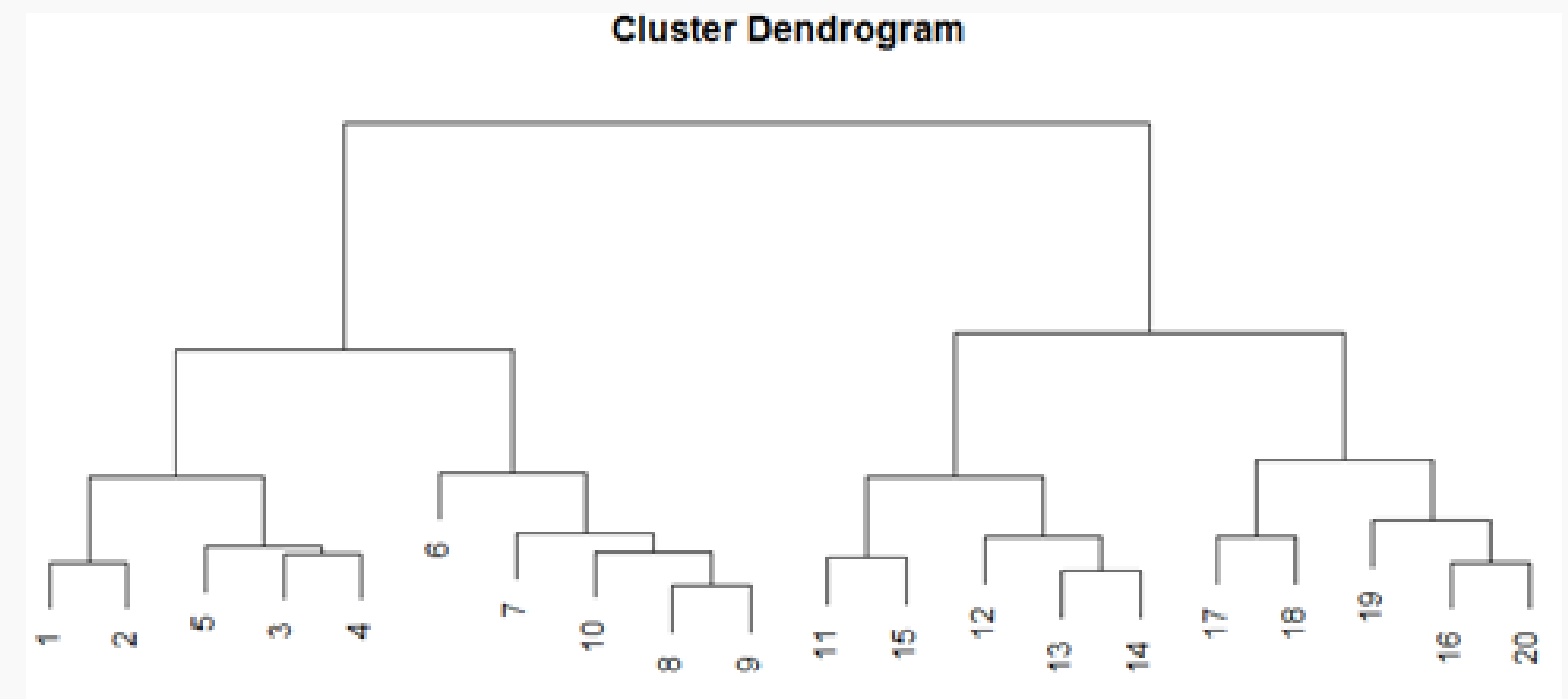
- 1、找到两个最近的两个点（类）；
- 2、把他们合成一个点（类）；（这个新的点不是数据集中本来存在的点，把原来那两个点删掉，用这个点代替原来两个点）
- 3、重复1、2。最终得到一棵树。



## 类间距离的计算

- 最短距离法
- 最长距离法
- 中间距离法
- 类平均法

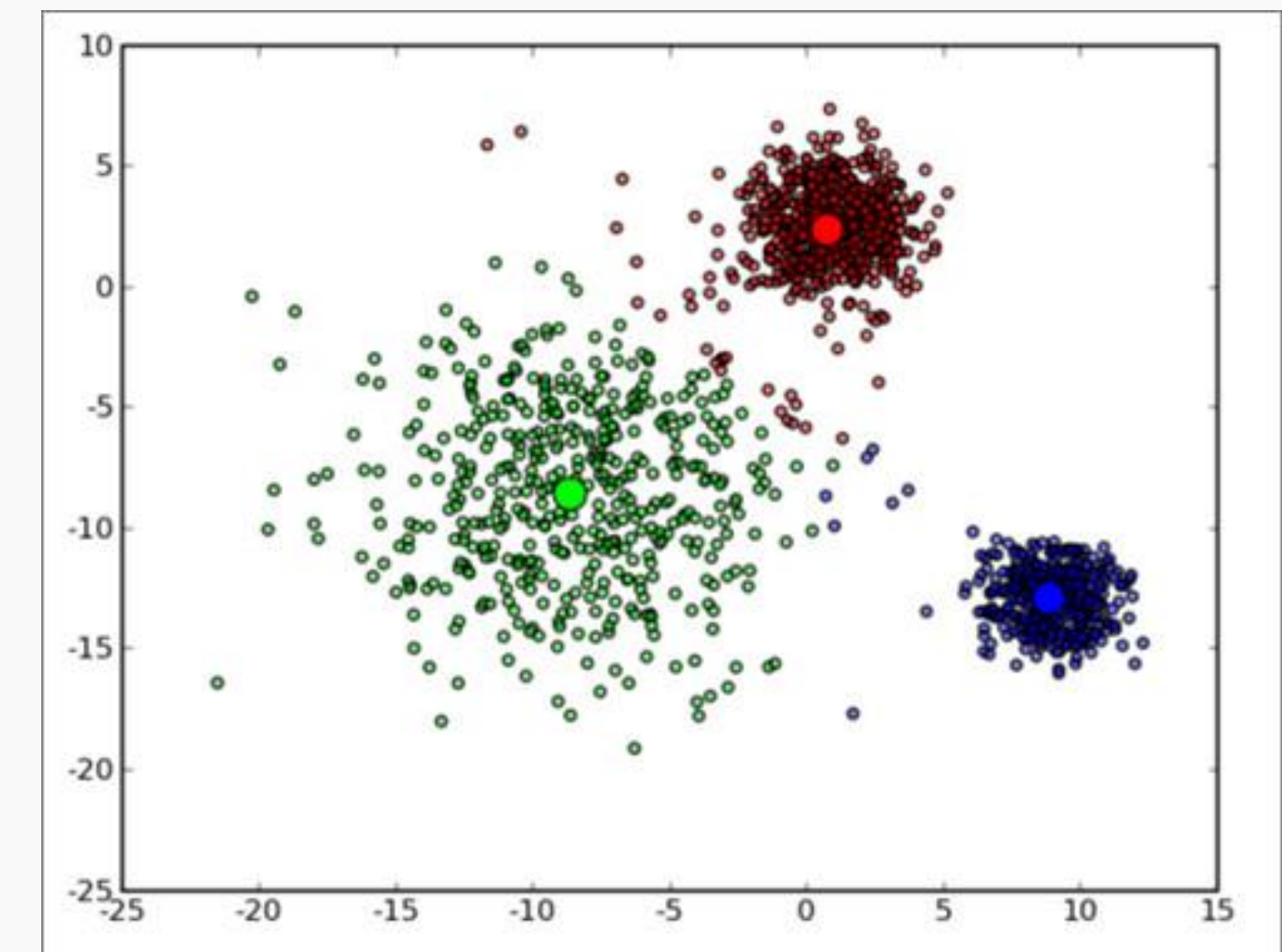
不同的距离方法，得到不同的结果





## K-means聚类法

- 1、随机选择k个点作为初始质心
- 2、把每个点按照距离分配给最近的质心，形成k个簇
- 3、重新计算每个簇的质心
- 4、重复2、3步，直到质心不再变化



## 数据变换

中心化变换：均值为0，方差不变

$$x_{ij}^* = x_{ij} - \bar{x}_j,$$

标准化变换：均值为0，方差为1

$$x_{ij}^* = \frac{x_{ij} - \bar{x}_j}{s_j}$$

极差正规化变换：极差为1，数值在[0,1]之间

$$x_{ij}^* = \frac{x_{ij} - \min_{1 \leq k \leq n} x_{kj}}{R_j}$$

# 实战

# 图像按照色彩聚类



# 聚类算法





# 聚类算法

