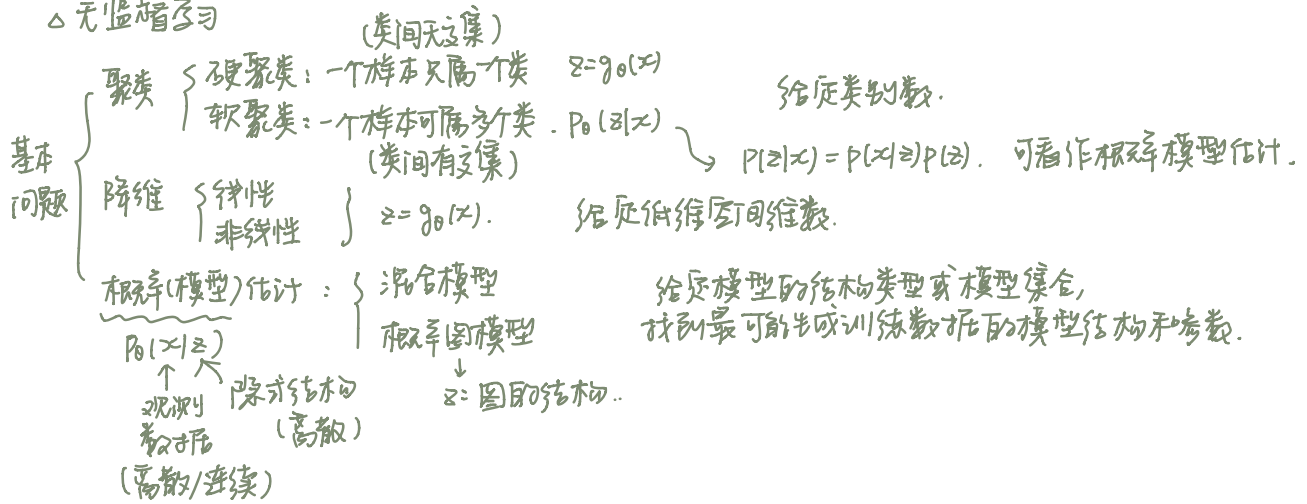


△ 无监督学习



△ 聚类

- 聚合(自下而上)聚类: 开始每个样本各属一类, 逐渐合并最近两类
- 分裂(自上而下)聚类: 开始所有样本同属一类, 逐渐将相距最远的样本分到两个新类.
- k-means 聚类: 基于中心.



△ 层次聚类 - 聚合聚类

- ① 计算 n 个样本两两间距离
- ② 构造 n 类, 每类只有一个样本
- ③ 合并类间距离最小的两类为一新类
- ④ 重复③ 至类数为 1.

△ k-means 聚类

划分方法为 C , $C(x_i)$ 表示 x_i 所在类.

target: $C^* = \arg \min_C \sum_{i=1}^K \sum_{x_i \in C_i} \|x_i - \bar{x}_i\|^2 = W(C)$ 能量. 表示相同类中样本相似程度. (越相似, 损失/能量越小) 1

用欧氏距离平方作样本间距离.

迭代求解：① 给定 k 个中心 (m_1, m_2, \dots, m_k) ，求划分 \hat{C} ：

$$\hat{C} = \arg \min_C \sum_{i=1}^k \sum_{x_j \in C_i} \|x_j - m_i\|^2$$

得： \hat{C} ：将每个样本指派到其最近的中心所在类。

② 给定划分 C ，再求各个类中心 $(\hat{m}_1, \dots, \hat{m}_k)$ ：

$$\hat{m}_1, \dots, \hat{m}_k = \arg \min_{m_1, \dots, m_k} \sum_{i=1}^k \sum_{x_j \in C_i} \|x_j - m_i\|^2$$

得： $\hat{m}_i = \frac{1}{n_i} \sum_{x_j \in C_i} x_j$ 类中心为该类样本均值。

步骤聚：① 随机选 k 个样本点作为初始中心。

② 分离样本至最近中心

③ 计算新中心（均值）

④ 收敛/符合停止条件，结束。否则返回②。
(划分不再改变)

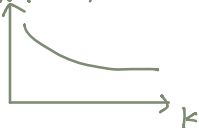
不保证收敛到全局最优。

初始中心的选取影响结果：可用层次聚类先对样本进行聚类，得 k 个类时停止，在哪个类选择

一个占其中 n_i (均值) 最接近的样本点。

最优 k ：实验。

类的平均直径。



平均直径越小，聚类结果质量越高。

二分查找找最优 k 值。