

谱聚类

聚类问题可以分为两种思路：

1. Compactness, 这类有 K-means, GMM 等, 但是这类算法只能处理凸集, 为了处理非凸的样本集, 必须引入核技巧。
2. Connectivity, 这类以谱聚类为代表。

谱聚类是一种基于无向带权图的聚类方法。这个图用 $G = (V, E)$ 表示, 其中 $V = \{1, 2, \dots, N\}$, $E = \{w_{ij}\}$, 这里 w_{ij} 就是边的权重, 这里权重取为相似度, $W = (w_{ij})$ 是相似度矩阵, 定义相似度 (径向核) :

$$\begin{aligned} w_{ij} = k(x_i, x_j) &= \exp\left(-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}\right), (i, j) \in E \\ w_{ij} &= 0, (i, j) \notin E \end{aligned} \quad (1)$$

下面定义图的分割, 这种分割就相当于聚类的结果。定义 $w(A, B)$:

$$A \subset V, B \subset V, A \cap B = \emptyset, w(A, B) = \sum_{i \in A, j \in B} w_{ij} \quad (2)$$

假设一共有 K 个类别, 对这个图的分割

$$CUT(V) = CUT(A_1, A_2, \dots, A_K) = \sum_{k=1}^K w(A_k, \overline{A_k}) = \sum_{k=1}^K [w(A_k, V) - w(A_k, A_k)]$$

于是, 我们的目标就是 $\min_{A_k} CUT(V)$ 。

为了平衡每一类内部的权重不同, 我们做归一化的操作, 定义每一个集合的度, 首先, 对单个节点的度定义:

$$d_i = \sum_{j=1}^N w_{ij} \quad (3)$$

其次, 每个集合:

$$\Delta_k = degree(A_k) = \sum_{i \in A_k} d_i \quad (4)$$

于是:

$$N(CUT) = \sum_{k=1}^K \frac{w(A_k, \overline{A_k})}{\sum_{i \in A_k} d_i} \quad (5)$$

所以目标函数就是最小化这个式子。

谱聚类的模型就是:

$$\{\hat{A}_k\}_{k=1}^K = \underset{A_k}{argmin} N(CUT) \quad (6)$$

引入指示向量：

$$\begin{cases} y_i \in \{0, 1\}^K \\ \sum_{j=1}^K y_{ij} = 1 \end{cases} \quad (7)$$

其中， y_{ij} 表示第 i 个样本属于 j 个类别，记： $Y = (y_1, y_2, \dots, y_N)^T$ 。所以：

$$\hat{Y} = \underset{Y}{\operatorname{argmin}} N(CUT) \quad (9)$$

将 $N(CUT)$ 写成对角矩阵的形式，于是：

$$\begin{aligned} N(CUT) &= \operatorname{Trace}[\operatorname{diag}(\frac{w(A_1, \overline{A_1})}{\sum_{i \in A_1} d_i}, \frac{w(A_2, \overline{A_2})}{\sum_{i \in A_2} d_i}, \dots, \frac{w(A_K, \overline{A_K})}{\sum_{i \in A_K} d_i})] \\ &= \operatorname{Trace}[\operatorname{diag}(w(A_1, \overline{A_1}), w(A_2, \overline{A_2}), \dots, w(A_K, \overline{A_K})) \cdot \operatorname{diag}(\sum_{i \in A_1} d_i, \dots, \sum_{i \in A_K} d_i)^{-1}] \\ &= \operatorname{Trace}[O \cdot P^{-1}] \end{aligned} \quad (10)$$

我们已经知道 Y, w 这两个矩阵，我们希望求得 O, P 。

由于：

$$Y^T Y = \sum_{i=1}^N y_i y_i^T \quad (11)$$

对于 $y_i y_i^T$ ，只在对角线上的 $k \times k$ 处为 1，所以：

$$Y^T Y = \operatorname{diag}(N_1, N_2, \dots, N_K) \quad (12)$$

其中， N_i 表示有 N_i 个样本属于 i ，即 $N_k = \sum_{k \in A_k} 1$ 。

引入对角矩阵，根据 d_i 的定义， $D = \operatorname{diag}(d_1, d_2, \dots, d_N) = \operatorname{diag}(w_{NN} \mathbb{I}_{N1})$ ，于是：

$$P = Y^T D Y \quad (13)$$

对另一项 $O = \operatorname{diag}(w(A_1, \overline{A_1}), w(A_2, \overline{A_2}), \dots, w(A_K, \overline{A_K}))$ ：

$$O = \operatorname{diag}(w(A_i, V)) - \operatorname{diag}(w(A_i, A_i)) = \operatorname{diag}(\sum_{j \in A_i} d_j) - \operatorname{diag}(w(A_i, A_i)) \quad (14)$$

其中，第一项已知，第二项可以写成 $Y^T w Y$ ，这是由于：

$$Y^T w Y = \sum_{i=1}^N \sum_{j=1}^N y_i y_j^T w_{ij} \quad (15)$$

于是这个矩阵的第 lm 项可以写为：

$$\sum_{i \in A_l, j \in A_m} w_{ij} \quad (16)$$

这个矩阵的对角线上的项和 $w(A_i, A_i)$ 相同，所以取迹后的取值不会变化。

所以：

$$N(CUT) = Trace[(Y^T(D - w))Y] \cdot (Y^T D Y)^{-1} \quad (17)$$

其中， $L = D - w$ 叫做拉普拉斯矩阵。