

# 谱聚类教程

乌尔里克·冯·卢克斯堡

马克斯普朗克生物控制论研究所

斯佩曼斯特 38, 72076 图宾根, 德国

ulrike.luxburg@tuebingen.mpg.de

本文发表于 Statistics and Computing, 17 (4), 2007。

原始出版物可在 [www.springer.com](http://www.springer.com)。

## 抽象的

近年来, 谱聚类已成为最流行的现代聚类算法之一。它实现简单, 可以通过标准线性代数软件有效地解决, 并且通常优于传统的聚类算法, 如  $k$ -means 算法。乍看之下, 谱聚类显得有些神秘, 并且根本看不出它为什么起作用以及它到底做了什么。本教程的目标是对这些问题给出一些直觉。我们描述了不同的图拉普拉斯算子及其基本属性, 介绍了最常见的谱聚类算法, 并通过几种不同的方法从头开始推导这些算法。讨论了不同谱聚类算法的优缺点。

关键词: 谱聚类; 图拉普拉斯算子

## 1 简介

聚类是探索性数据分析中使用最广泛的技术之一, 应用范围从统计学、计算机科学、生物学到社会科学或心理学。实际上, 在处理经验数据的每个科学领域中, 人们都试图通过识别数据中的“相似行为”组来获得对数据的第一印象。在本文中, 我们想向读者介绍谱聚类算法系列。与“传统算法”相比, 例如  $k$ -手段或单链接, 谱聚类具有许多基本优势。通过谱聚类获得的结果通常优于传统方法, 谱聚类实现起来非常简单, 并且可以通过标准线性代数方法有效地求解。

本教程是对谱聚类的独立介绍。我们从头开始推导出谱聚类, 并就谱聚类为何起作用提出了不同的观点。除了基本的线性代数, 读者不需要特定的数学背景。然而, 我们并没有试图对所有关于谱聚类的文献进行简明的回顾, 这是不可能的, 因为关于这个主题的文献太多了。前两节专门逐步介绍谱聚类使用的数学对象: 第 2 节中的相似图和第 3 节中的拉普拉斯图。谱聚类算法本身将在第 4 节中介绍。下一个然后用三个部分专门解释为什么这些算法有效。每个部分对应一个解释: 第 5 部分描述了一种图形划分方法, 第 6 部分描述了随机游走的观点, 第 7 部分描述了扰动理论方法。在第 8 节中, 我们将研究与谱聚类相关的一些实际问题, 并在第 9 节中讨论与谱聚类相关的各种扩展和文献。

## 2 相似度图

给定一组数据点  $x_1, \dots, x_n$  和一些相似性的概念  $s_{ij}$ ，所有数据点对之间为  $0 \leq s_{ij} \leq 1$  和  $x_i, x_j$ ，聚类的直观目标是将数据点分成几组，使得同一组中的点彼此相似，不同组中的点彼此不相似。如果我们没有比数据点之间的相似性更多的信息，表示数据的一种很好的方式是 **相似图**  $G = (V, E)$ 。每个顶点  $v_i$  在此图中代表一个数据点  $x_i$ 。如果相似性则两个顶点相连  $s_{ij}$  对应数据点之间  $x_i$  和  $x_j$  为正值或大于某个阈值，边缘被加权  $s_{ij}$ 。现在可以使用相似度图重新表述聚类问题：我们想要找到图的一个分区，使得不同组之间的边具有非常低的权重（这意味着不同聚类中的点彼此不同）并且边组内的点具有高权重（这意味着同一簇内的点彼此相似）。为了能够形式化这种直觉，我们首先要介绍一些基本的图形符号，并简要讨论我们将要研究的图形类型。

### 2.1 图形符号

让  $G = (V, E)$  是一个有顶点集的无向图  $V = \{v_1, \dots, v_n\}$ 。下面我们假设图  $G$  是加权的，即两个顶点之间的每条边  $v_i$  和  $v_j$  具有非负权重  $w_{ij} \geq 0$ 。加权 **邻接矩阵图** 的是矩阵  $W = (w_{ij})_{i,j=1,\dots,n}$ 。如果  $w_{ij} = 0$  这意味着顶点  $v_i$  和  $v_j$  没有边连接。作为  $G$  是无向的，我们需要  $w_{ij} = w_{ji}$ 。顶点的度数  $d_i \in V$  定义为

$$d_i = \sum_{j=1}^n w_{ij}.$$

注意，事实上，这个和只遍历所有相邻的顶点  $v_j$ ，对于所有其他顶点  $v_j$  重量  $w_{ij}$  为 0。**度矩阵**  $D$  被定义为具有度数的对角矩阵  $d_1, \dots, d_n$  在对角线上。给定一个顶点子集  $A \subset V$ ，我们表示它的补码  $V \setminus A$  经过一个。我们定义指标向量  $\mathbf{1}_A = (f_1, \dots, f_n) \in \mathbb{R}^n$  作为带有条目的向量  $f_i = 1$  如果  $v_i \in A$  和  $f_i = 0$  否则。为了方便我们引入  $\text{sh}$

手写符号  $\mathbf{1}_A \in \mathbb{R}^n$  对于一组索引  $A \subset V$ 。对于两个不一定不相交的集合  $A, B \subset V$  我们定义

$\{v_i \mid v_i \in A\}$ ，特别是在处理像  $A, B \subset V$  我们定义

$$W(A, B) := \sum_{i \in A, j \in B} w_{ij}.$$

我们考虑两种不同的方法来衡量子集的“大小”  $A \subset V$ ：

$$\begin{aligned} |A| &:= \text{中的顶点数 } A \\ \text{体积}(A) &:= \sum_{i \in A} d_i. \end{aligned}$$

直觉上， $|A|$  测量的大小  $A$  通过它的顶点数，而  $\text{vol}(A)$  测量的大小  $A$  通过对附加到顶点的所有边的权重求和  $A$ 。一个子集  $A \subset V$  如果图中的任意两个顶点是连通的  $A$  可以通过路径连接，使得所有中间点也位于  $A$ 。一个子集  $A$  如果它是连通的并且如果在  $A$  中的顶点之间没有连接，则称为连通分量  $A$  和一个。非空集  $A_1, \dots, A_k$  形成图的分区，如果  $A_i \cap A_j = \emptyset$  和  $A_1 \cup \dots \cup A_k = V$ 。

## 2.2 不同的相似度图

有几种流行的构造来转换给定的集合  $x_1, \dots, x_n$  具有成对相似性的数据点  $s_{ij}$  或成对距离  $d_{ij}$  成图表。构建相似度图时，目标是对数据点之间的局部邻域关系进行建模。

**$\epsilon$ -邻域图**：这里我们连接成对距离小于  $\epsilon$  的所有点。由于所有连接点之间的距离大致具有相同的比例（至多  $\epsilon$ ），对边缘进行加权不会将有关数据的更多信息合并到图中。因此， $\epsilon$ -邻域图通常被认为是未加权的图。

**$k$ -最近邻图**：这里的目标是连接顶点  $v_i$  和顶点  $v_j$  如果  $v_j$  是其中之一  $k$ -最近的邻居  $v_i$ 。然而，这个定义导致有向图，因为邻域关系不是对称的。有两种方法可以使这个图无向。第一种方法是简单地忽略边缘的方向，即我们连接  $v_i$  和  $v_j$  如果有无向边  $v_i$

是其中之一  $k$ -最近的邻居  $v_j$  或者如果  $v_j$  是其中之一  $k$ -最近的邻居  $v_i$ 。结果图就是通常所说的这  $k$ -最近邻图。第二种选择是连接顶点  $v_i$  和  $v_j$  如果两者  $v_i$  是其中之一  $k$ -最近的邻居  $v_j$  和  $v_j$  是其中之一  $k$ -最近的邻居  $v_i$ 。生成的图形称为 *相互的*  $k$ -最近邻图。在这两种情况下，在连接适当的顶点后，我们根据端点的相似性对边进行加权。

**全连接图**：在这里，我们简单地将所有具有正相似性的点彼此连接起来，然后我们对所有边进行加权  $s_{ij}$ 。由于该图应该表示局部邻域关系，因此这种构造仅在相似性函数本身对局部邻域建模时才有用。这种相似性函数的一个例子是高斯相似性函数  $s(x_i, x_j) = \exp(-\|x_i - x_j\|_2^2 / (2\sigma^2))$ ，其中参数  $\sigma$  控制邻域的宽度。该参数起着与参数类似的作用  $\epsilon$  在的情况下  $\epsilon$ -邻域图。

上面提到的所有图都经常用于谱聚类。据我们所知，关于相似图的选择如何影响谱聚类结果的问题的理论结果并不存在。关于不同图的行为的讨论，我们参考第 8 节。

## 3 图拉普拉斯算子及其基本性质

谱聚类的主要工具是图拉普拉斯矩阵。存在一个致力于研究这些矩阵的整个领域，称为谱图理论（例如，参见 Chung, 1997）。在本节中，我们要定义不同的图拉普拉斯算子并指出它们最重要的属性。我们将仔细区分图拉普拉斯算子的不同变体。请注意，在文献中没有唯一的约定将哪个矩阵确切地称为“拉普拉斯图”。通常，每个作者只是将“他的”矩阵称为拉普拉斯图。因此，在阅读有关图拉普拉斯算子的文献时需要格外小心。

在下文中，我们总是假设  $G$  是具有权重矩阵的无向加权图  $W$ ，在哪里  $w_{ij} = w_{ji} \geq 0$ 。当使用矩阵的特征向量时，我们不一定假设它们是归一化的。例如，常数向量  $\mathbf{1}$  和多个  $\mathbf{1}$  对于一些  $c \neq 0$  将被视为相同的特征向量。特征值将始终递增排序，尊重多重性。通过“第  $k$  特征向量”，我们指的是对应于  $k$  最小的特征值。

### 3.1 非归一化图拉普拉斯算子

非归一化图拉普拉斯矩阵定义为

$$L = D - W.$$

在 Mohar (1991, 1997) 中可以找到对其许多属性的概述。以下命题总结了谱聚类所需的最重要的事实。

提案 1 (属性  $L$ ) 矩阵  $L$  满足以下属性:

1. 对于每个向量  $F \in \mathbb{R}^n$  我们有

$$F^T L F = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (F_i - F_j)^2.$$

2.  $L$  是对称的和半正定的。

3.  $L$  的最小特征值为 0, 对应的特征向量为常数一向量  $\mathbf{1}$ 。

4.  $L$  有  $n$  非负实值特征值  $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ 。

证明。

第 (1) 部分: 根据定义  $L = D - W$ ,

$$\begin{aligned} F^T L F &= F^T (D - W) F = \sum_{i=1}^n d_i F_i^2 - \sum_{i,j=1}^n F_i F_j w_{ij} \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (d_i + d_j - 2w_{ij}) F_i F_j = \frac{1}{2} \sum_{i,j=1}^n (d_i + d_j - 2w_{ij}) F_i F_j \\ &= \frac{1}{2} \sum_{i,j=1}^n (d_i + d_j - 2w_{ij}) F_i F_j = \frac{1}{2} \sum_{i,j=1}^n (d_i + d_j - 2w_{ij}) F_i F_j \end{aligned}$$

第 (2) 部分: 的对称性  $L$  直接从对称性得出  $W$  和  $D$ . 半正定性是第 (1) 部分的直接结果, 它表明  $F^T L F \geq 0$  对于所有  $F \in \mathbb{R}^n$ . 第 (3) 部分: 显而易见。

第 (4) 部分是第 (1) - (3) 部分的直接结果。

2个

请注意, 未归一化图拉普拉斯算子不依赖于邻接矩阵的对角线元素  $W$ . 每个邻接矩阵与  $W$  在所有非对角线位置上导致相同的非标准化图拉普拉斯算子  $L$ . 特别是, 图中的自边不会改变相应的图拉普拉斯算子。

非标准化图拉普拉斯算子及其特征值和特征向量可用于描述图的许多属性, 参见 Mohar (1991, 1997)。下面是一个对谱聚类很重要的例子:

命题 2 (连通分量的个数和  $L$ ) 让  $G$  是具有非负权重的无向图。然后  $L$  具有多重性  $k$  的特征值 0 的个数等于连接组件的数量  $k_1, \dots, k_k$  在图中。特征值 0 的特征空间由指标向量跨越  $\mathbf{1}_1, \dots, \mathbf{1}_k$  这些组件。

$k$

证明。我们从案例入手  $k=1$ , 即图是连通的。假使, 假设  $F$  是一个特征值为 0 的特征向量。那么我们知道

$$0 = F^T L F = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (F_i - F_j)^2.$$

4个

作为权重  $w_{ij}$  是非负的，这个总和只能消失，如果所有条款  $w_{ij}(F_i - F_j)$  消失。因此，如果两个顶点  $v_i$  和  $v_j$  连接（即， $w_{ij} > 0$ ），那么  $F_i$  需要等于  $F_j$ 。有了这个论点，我们可以看到  $F$  对于可以通过图中的路径连接的所有顶点，需要保持不变。此外，由于无向图中连通分量的所有顶点都可以通过路径连接， $F$  需要在整个连接的组件上保持不变。在仅由一个连通分量组成的图中，我们因此只有一个常量向量  $\mathbf{1}$  为特征值为 0 的特征向量，显然是连通分量的指示向量。

现在考虑的情况  $k$  连接的组件。在不失一般性的情况下，我们假设顶点是根据它们所属的连通分量排序的。在这种情况下，邻接矩阵  $W$  有分块对角线形式，矩阵也一样  $L$ ：

$$L = \begin{pmatrix} L_1 & & \\ & L_2 & \\ & & \ddots \\ & & & L_k \end{pmatrix}$$

请注意，每个块  $L_i$  本身是一个适当的图拉普拉斯算子，即对应于的子图的拉普拉斯算子  $i$ -th 连接组件。对于所有块对角矩阵都是这种情况，我们知道  $L$  由光谱的并集给出  $L_i$ ，以及相应的特征向量  $L_i$  是特征向量  $L_i$ ，在其他块的位置填充 0。作为每个  $L_i$  是连通图的拉普拉斯图，我们知道每个  $L_i$  具有重数为 1 的特征值 0，对应的特征向量是  $i$ -th 连接组件。因此，矩阵  $L$  具有与连通分量一样多的特征值 0，对应的特征向量是连通分量的指示向量。

2个

### 3.2 归一化图拉普拉斯算子

文献中有两个矩阵称为归一化图拉普拉斯矩阵。这两个矩阵彼此密切相关，定义为

$$L_{\text{sym}} := \frac{1}{2} (D^{-1/2} L D^{-1/2}) = \frac{1}{2} (D - W) D^{-1/2}$$

$$L_{\text{rw}} := \frac{1}{n} D^{-1} L$$

我们将第一个矩阵表示为  $L_{\text{sym}}$  因为它是一个对称矩阵，第二个是  $L_{\text{rw}}$  因为它与随机游走密切相关。下面我们总结了几个属性  $L_{\text{sym}}$  和  $L_{\text{rw}}$ 。规范化图拉普拉斯算子的标准参考文献是 Chung (1997)。

提案 3 (的属性  $L_{\text{sym}}$  和  $L_{\text{rw}}$ ) 归一化拉普拉斯算子满足以下属性：

1. 对于每个  $F \in \mathbb{R}^n$  我们有

$$F^T L_{\text{sym}} F = \frac{1}{2} \sum_{i,j=1}^n w_{ij} \left( \frac{F_i}{\sqrt{d_i}} - \frac{F_j}{\sqrt{d_j}} \right)^2$$

2.  $\lambda$  是  $L_{\text{sym}}$  的特征值与特征向量  $w$  当且仅当  $\lambda$  是  $L_{\text{rw}}$  与特征向量  $w = \frac{1}{\sqrt{2}} \mathbf{1}$ 。

3.  $\lambda$  是一个特征值  $L_{\text{rw}}$  与特征向量  $w$  当且仅当  $\lambda$  和你解决广义特征问题  $Lw = \lambda Dw$ 。

4.0 是特征值  $\lambda_1$  与常数  $\frac{1}{n}$  作为特征向量。0 是特征值  $\lambda_2$  与特征向量  $\frac{1}{n}$ 。

5. 大符号  $L$  和大号  $L$  是半正定的并且有  $n$  非负实值特征值  $0 = \lambda_1 \leq \dots \leq \lambda_n$ 。

证明。第 (1) 部分可以类似于命题 1 的第 (1) 部分来证明。

乘以特征值方程可以立即看到部分 (2) 大符号  $w = \lambda w$  和  $\frac{1}{n}$ 。

从左边开始代入  $w = \frac{1}{n}$ 。

第 (3) 部分直接乘以特征值方程  $w = \lambda w$  和  $\frac{1}{n}$  从左边。第 (4) 部分：第一个陈述是显而易见的大号  $L$  字  $1 = 0$ ，第二个陈述来自 (2)。第 (5) 部分：关于声明大符号从 (1) 得出，然后是关于大号  $L$  从 (2) 得出。

2个

与非归一化图拉普拉斯算子的情况一样，归一化图拉普拉斯算子的特征值 0 的重数与连通分量的数量有关：

命题 4 (的连通分量数和光谱大符号和大号  $L$ ) 让  $G$  是具有非负权重的无向图。然后是多重性  $k$  的特征值 0 两者的大号  $L$  字

和大号  $L$  等于连接组件的数量  $1, \dots, k$  在图中。为了大号  $L$  字, 本征空间 0 由指标向量跨越  $1, \dots, k$  这些组件。为了大符号, 本征空间 0 被向量跨越  $\frac{1}{n}, \dots, \frac{1}{n}$ 。

证明。证明类似于命题 2 的证明，使用命题 3。

2个

## 4 谱聚类算法

现在我们想说明最常见的谱聚类算法。对于参考和光谱聚类的历史，我们参考第 9 节。我们假设我们的数据包括  $n$  “点”  $x_1, \dots, x_n$  可以是任意对象。我们测量它们的成对相似性  $s_{ij} = s(x_i, x_j)$  通过一些对称且非负的相似性函数，我们将相应的相似性矩阵表示为  $S = (s_{ij})_{i,j=1,\dots,n}$ 。

### 非归一化谱聚类

输入：相似度矩阵  $S \in \mathbb{R}^{n \times n}$ ，数字  $k$  要构建的集群。

- 通过第 2 节中描述的方法之一构建相似图。作为其加权邻接矩阵。

让  $W$

- 计算非标准化拉普拉斯算子  $L$ 。
- 计算第一个  $k$  特征向量  $w_1, \dots, w_k$  的大号。
- 让  $U \in \mathbb{R}^{n \times k}$  是包含向量的矩阵  $w_1, \dots, w_k$  作为列。
- 为了一  $i = 1, \dots, n$ ，让  $z_i \in \mathbb{R}^k$  是对应于  $i$  第行  $U$ 。
- 聚焦点  $(z_1, \dots, z_n)$  在  $\mathbb{R}^k$  与  $k$ -means 算法成簇  $C_1, \dots, C_k$ 。

输出：集群  $1, \dots, k$  和一个  $i = \{j \mid z_j \in C_i\}$ 。

归一化谱聚类有两种不同的版本，具体取决于哪个归一化的

使用图拉普拉斯算子。我们以两篇热门论文命名这两种算法，更多参考资料和历史请参阅第 9 节。

根据 Shi 和 Malik (2000) 的归一化光谱聚类

输入：相似度矩阵  $S \in \mathbb{R}^{n \times n}$ ，数字  $k$  要构建的集群。

- 通过第 2 节中描述的方法之一构建相似图。作为其加权邻接矩阵。

让  $W$

- 计算非标准化拉普拉斯算子  $L$ 。
- 计算第一个  $k$  广义特征向量  $u_1, \dots, u_k$  广义特征问题  $Lx = \lambda x$ 。
- 让  $\tilde{U} \in \mathbb{R}^{n \times k}$  是包含向量的矩阵  $u_1, \dots, u_k$  作为列。
- 为了一世  $i=1, \dots, \tilde{n}$ ，让  $\tilde{c}_i \in \mathbb{R}^k$  是对应于一世-第行  $\tilde{u}_i$ 。
- 聚集点  $(\tilde{c}_1, \dots, \tilde{c}_{\tilde{n}})$  在  $\mathbb{R}^k$  与  $k$ -means 算法成簇  $C_1, \dots, C_k$ 。

输出：集群一个  $1, \dots, k$  和一个  $\tilde{c}_i = \{j \mid j \in C_{\tilde{c}_i}\}$ 。

请注意，此算法使用的广义特征向量  $u_i$ ，根据命题 3 对应于矩阵的特征向量  $u_i$  写字。所以实际上，该算法适用于归一化拉普拉斯算子的特征向量  $u_i$  写字，因此称为归一化谱聚类。下一个算法也使用归一化的拉普拉斯算子，但这次矩阵  $S$  符号代替  $S$  写字。正如我们将看到的，该算法需要引入其他算法不需要的额外行规范化步骤。原因将在第 7 节中变得清晰。

根据 Ng、Jordan 和 Weiss (2002) 的规范化光谱聚类

输入：相似度矩阵  $S \in \mathbb{R}^{n \times n}$ ，数字  $k$  要构建的集群。

- 通过第 2 节中描述的方法之一构建相似图。作为其加权邻接矩阵。

让  $W$

- 计算归一化拉普拉斯算子  $L$  符号。
- 计算第一个  $k$  特征向量  $u_1, \dots, u_k$  的大号符号。
- 让  $\tilde{U} \in \mathbb{R}^{n \times k}$  是包含向量的矩阵  $u_1, \dots, u_k$  作为列。
- 形成矩阵  $\tilde{L} \in \mathbb{R}^{n \times k}$   $\sum_{i=1}^k \tilde{L}_{ij} = 1$  从  $\tilde{u}_i$  通过将行规范化为范数 1，那是设定的  $\tilde{L}_{ij} = \tilde{u}_{ij} / \sqrt{\sum_{l=1}^k \tilde{u}_{il}^2}$  我知道 1 个 2 个。
- 为了一世  $i=1, \dots, \tilde{n}$ ，让  $\tilde{c}_i \in \mathbb{R}^k$  是对应于一世-第行  $\tilde{L}_{i\cdot}$ 。
- 聚集点  $(\tilde{c}_1, \dots, \tilde{c}_{\tilde{n}})$  与  $k$ -means 算法成簇  $C_1, \dots, C_k$ 。输出：集群一个  $1, \dots, k$  和一个  $\tilde{c}_i = \{j \mid j \in C_{\tilde{c}_i}\}$ 。

除了使用三种不同的图拉普拉斯算子之外，上述所有三种算法看起来都非常相似。在这三种算法中，主要技巧是改变抽象数据点的表示  $X_{\tilde{c}_i}$  到点  $\tilde{c}_i \in \mathbb{R}^k$ 。由于图拉普拉斯算子的特性，这种表示形式的变化很有用。我们将在下一节中看到这种表示形式的变化增强了数据中的聚类属性，因此可以在新表示形式中轻松检测到聚类。特别是，简单的  $k$ -意味着聚类算法可以毫无困难地检测这种新表示中的聚类。不熟悉的读者  $k$ -意味着可以在很多地方阅读这个算法

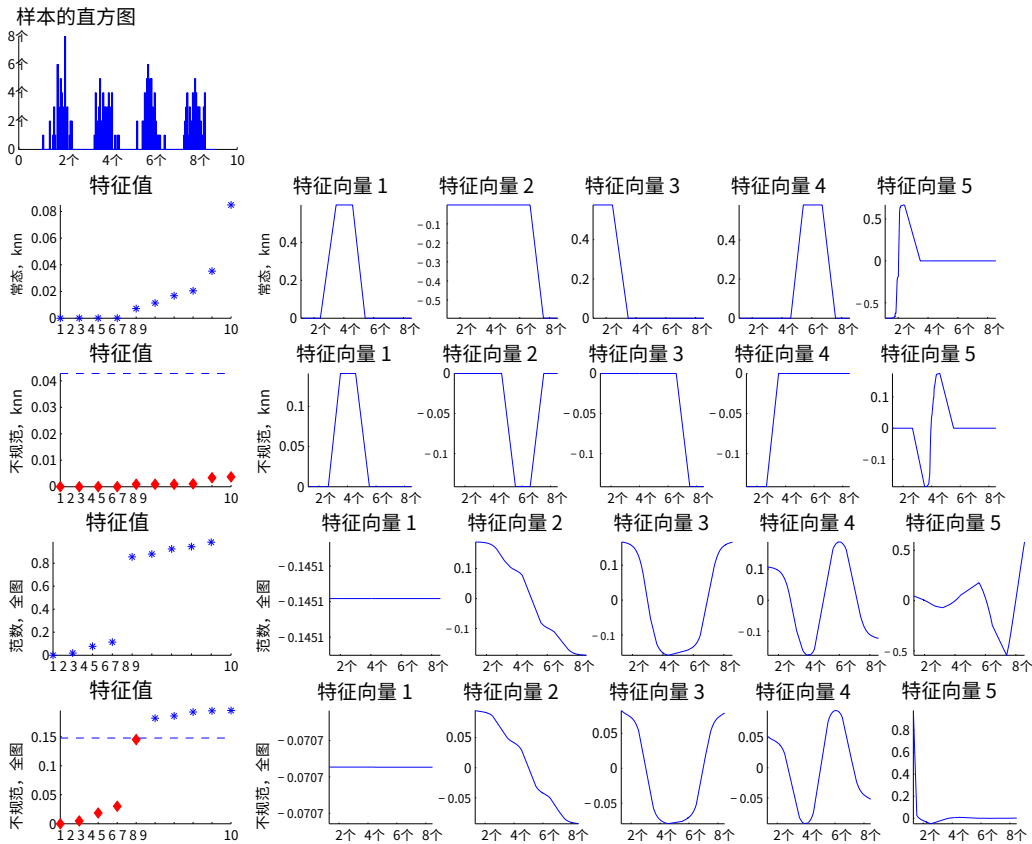


图 1: 光谱聚类的玩具示例, 其中数据点是从四个高斯分布的混合中提取的R。左上角: 数据的直方图。第一行和第二行: 的特征值和特征向量大写字母和基于k-最近邻图。第三行和第四行: 的特征值和特征向量大写字母和基于全连接图。对于所有图, 我们使用高斯核 $\sigma = 1$ 作为相似函数。有关详细信息, 请参阅文本。

教科书, 例如 Hastie、Tibshirani 和 Friedman (2001)。

在我们深入研究谱聚类理论之前, 我们想通过一个非常简单的玩具示例来说明其原理。这个例子将在本教程的几个地方使用, 我们选择它是因为它非常简单, 相关量很容易绘制出来。这个玩具数据集由 200 个点的随机样本组成 $X_1, \dots, X_{200} \in \mathbb{R}$ 根据四个高斯的混合绘制。图 1 的第一行显示了从该分布中抽取的样本的直方图 ( $X$ -axis 表示一维数据空间)。作为该数据集的相似度函数, 我们选择高斯相似度函数 $\rho(X_i, X_j) = \exp(-\frac{1}{2\sigma^2}(X_i - X_j)^2)$ 和 $\sigma = 1$ 。全连接图和10-最近邻图。

作为相似图, 我们考虑了在图 1 中我们显示了第一个特征值

和非标准化拉普拉斯算子的特征向量大写字母和归一化的拉普拉斯算子大写字母。也就是说, 在我们绘制的特征值图中 $\lambda_i$ 对比 $\lambda_i$ (暂时忽略图中未归一化情况下的虚线和特征值的不同形状; 它们的含义将在第 8.5 节中讨论)。在特征向量的特征向量图中 $f = (f_1, \dots, f_{200})$ 我们密谋 $X_i$ 对比 $f_i$ (请注意, 在所选示例中 $X_i$ 只是一个实数, 因此我们可以在 $X$ -axis)。图 1 的前两行显示了基于 10 最近邻图的结果。我们可以看到前四个特征值

为 0, 并且 corresponding eigenvector 是 cluster 指标向量。原因是集群



在 10 最近邻图中形成断开连接的部分，在这种情况下，特征向量如命题 2 和 4 中给出。接下来的两行显示了完全连接图的结果。由于高斯相似函数始终为正，因此该图仅由一个连通分量组成。因此，特征值 0 的重数为 1，第一个特征向量是常数向量。以下特征向量携带有关集群的信息。例如，在未归一化的情况下（最后一行），如果我们将第二个特征向量阈值设置为 0，则 0 以下的部分对应于聚类 1 和 2，而 0 以上的部分对应于聚类 3 和 4。类似地，对第三个特征向量进行阈值分离集群 1 和 4 来自集群 2 和 3，并且对第四个特征向量进行阈值处理将集群 1 和 3 与集群 2 和 4 分开。总而言之，前四个特征向量包含有关四个簇的所有信息。在此图中说明的所有情况下，谱聚类使用  $k$ -means 在前四个特征向量上很容易检测到正确的四个簇。

## 5 图切点

聚类的直觉是根据点的相似性将点分成不同的组。对于以相似图形式给出的数据，这个问题可以重述如下：我们想要找到图的一个分区，使得不同组之间的边具有非常低的权重（这意味着不同集群中的点彼此）并且组内的边具有高权重（这意味着同一簇内的点彼此相似）。在本节中，我们将看到如何将谱聚类导出为此类图划分问题的近似值。

给定一个具有邻接矩阵的相似图  $W$ ，最简单直接的构造分区的方法  $\Sigma$  该图是为了解决 mincut 问题。要定义它，请记住符号  $W(A, B) :=$

$\sum_{i \in A, j \in B} w_{ij}$  和一个为了补充一个  $k$  对于给定的数字  $k$  的子集，  
mincut 方法仅在于选择一个分区  $\{A_1, \dots, A_k\}$  最小化

$$\text{cut}(A_1, \dots, A_k) := \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k W(A_i, A_j) \mathbb{1}_{i \neq j}.$$

这里我们引入因子  $1/2$  为了符号的一致性，否则我们将在切割中计算每条边两次。特别是对于  $k=2$ ，mincut 是一个相对容易的问题，可以有效地解决，参见 Stoer 和 Wagner (1997) 及其中的讨论。然而，在实践中，它通常不会导致令人满意的分区。问题在于，在许多情况下，mincut 的解决方案只是将一个单独的顶点与图中的其余部分分开。当然这不是我们想要在聚类中实现的，因为聚类应该是相当大的点组。绕过这个问题的一种方法是明确要求集合  $A_1, \dots, A_k$  是“相当大”的。对此进行编码的两个最常见的目标函数是 RatioCut (Hagen 和 Kahng, 1992) 和归一化切割 Ncut (Shi 和 Malik, 2000)。在 RatioCut 中，子集的大小  $|A_i|$  的数量是通过它的顶点数来衡量的  $|A_i|$ ，而在 Ncut 中，尺寸是通过其边缘  $\text{vol}(A_i)$  定义的：

$$\begin{aligned} \text{比例切割}(A_1, \dots, A_k) &:= \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \frac{W(A_i, A_j)}{|A_i| |A_j|} = \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{|A_i|} \\ \text{切割}(A_1, \dots, A_k) &:= \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \frac{W(A_i, A_j)}{\text{vol}(A_i) \text{vol}(A_j)} = \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{\text{vol}(A_i)}. \end{aligned}$$

请注意，这两个目标函数  $\Sigma$  如果集群取一个小值  $|A_i|$  不是太小。特别是  $\sum_{i=1}^k (1/|A_i|)$  如果全部实现  $|A_i|$  重合，并且最小的  $\sum_{i=1}^k (1/\text{vol}(A_i))$  如果所有  $\text{vol}(A_i)$  重合。所以这两个目标函数试图实现的是集群是“平衡的”，分别通过顶点数或边权重来衡量。不幸的是，引入平衡条件使得之前解决 mincut 问题变得简单

### 5.1 近似 RatioCut $k=2$ 个

$$\min_{\substack{\text{一个 } CV \\ \text{一个 } \overline{CV}}} \text{最小比率切割 } (\text{一个}, \text{一个}). \quad (1)$$
$$F_{-世} = \frac{-\sqrt{\frac{|\text{一个}||\text{一个}|}{|\text{一个}||\text{一个}|}}}{-\sqrt{\frac{|\text{一个}||\text{一个}|}{|\text{一个}||\text{一个}|}}} \quad \begin{array}{l} \text{如果 } v_{-世} \in \text{一个} \\ \text{如果 } v_{-世} \in \text{一个} \end{array} \quad (2)$$
$$\begin{aligned}
F \text{ 如果} &= \frac{1}{2} \sum_{i,j=1}^2 w_{ij} (F - \text{世} - F_j)^{2\uparrow} \\
&= \frac{1}{2} \sum_{-\text{世} \in j \in -\text{个}} w_{ij} \frac{\sqrt{F - \text{个}}}{\sqrt{F - \text{个}}} \frac{\sqrt{F - \text{个}}}{\sqrt{F - \text{个}}} + \frac{1}{2} \sum_{-\text{世} \in j \in -\text{个}} w_{ij} \frac{\sqrt{F - \text{个}}}{\sqrt{F - \text{个}}} \frac{\sqrt{F - \text{个}}}{\sqrt{F - \text{个}}} \\
&= \text{切} \left( \text{一个}, \text{一个} \right) \frac{\sqrt{F - \text{个}}}{\sqrt{F - \text{个}}} \frac{\sqrt{F - \text{个}}}{\sqrt{F - \text{个}}} \\
&= \text{切} \left( \text{一个}, \text{一个} \right) \frac{\sqrt{F - \text{个}}}{\sqrt{F - \text{个}}} + \frac{\sqrt{F - \text{个}}}{\sqrt{F - \text{个}}} \\
&= /5/ \cdot \text{比例切割} (\text{一个}, \text{一个}).
\end{aligned}$$
$$\sum_{-1 \leq j \leq 1} F_{-1-j} = \frac{\sum_{-1 \leq j \leq 1} \sqrt{-1-j}}{\sqrt{-1-j}} = \frac{\sum_{-1 \leq j \leq 1} \sqrt{-1-j}}{\sqrt{-1-j}} = \frac{\sqrt{-1-j}}{\sqrt{-1-j}} = 1.$$
$$\|F\|_{2\pi} = \sum_{\substack{F_2 \text{ 全 } \rightarrow \text{ 一个} \\ \text{一} \# = 1}} \frac{\overline{F_1} \uparrow}{F_1 \downarrow} / \frac{\overline{F_2} \uparrow}{F_2 \downarrow} / \frac{\overline{F_3} \uparrow}{F_3 \downarrow} / \frac{\overline{F_4} \uparrow}{F_4 \downarrow} / \text{一个} \downarrow / \text{一个} \downarrow \text{名词}$$

分钟  $F$  如果受制于  $F \perp 1$  个,  $F$  也如方程式中所定义。(2),  $\|F\| = \text{名词}$   $\sqrt{-}$  (3)

10

丢弃离散条件，而是允许  $F_{-i}$  取任意值  $R$ 。这导致松弛的优化问题

$$\min_{F \in \mathbb{R}^n} \|F\| \quad \text{如果受制于 } F \perp 1, \|F\| = \text{名词} \quad (4)$$

根据 Rayleigh-Ritz 定理（例如，参见 Lütkepohl, 1997 年的第 5.5.2 节），可以立即看出该问题的解由向量给出  $F$  这是对第二小特征值的特征向量（回想一下，最小的特征值为 0 与特征向量  $1$ ）。所以我们可以通过的第二个特征向量来近似 RatioCut 的最小值。然而，为了获得图的分区，我们需要重新转换实值解向量  $F$  将松弛问题转化为离散指标向量。最简单的方法是使用符号  $F$  作为指标函数，即选择

$$\begin{cases} v_{-i} \in \text{一个} & \text{如果 } F_{-i} \geq 0 \\ v_{-i} \in \text{一个} & \text{如果 } F_{-i} < 0. \end{cases}$$

但是，特别是在  $k > 2$  下面处理，这个启发式太简单了。大多数谱聚类算法所做的是考虑坐标  $F_{-i}$  作为点  $R$  并将他们分成两组  $C$ 。由  $k$ -means 聚类算法。然后将得到的聚类结果传递给底层数据点，也就是我们选择

$$\begin{cases} v_{-i} \in \text{一个} & \text{如果 } F_{-i} \in C \\ v_{-i} \in \text{一个} & \text{如果 } F_{-i} \in C^c. \end{cases}$$

这正是 **非归一化谱聚类算法**的情况  $k=2$ 。

## 5.2 任意近似 RatioCut $k$

一般值情况下 RatioCut 最小化问题的松弛  $k$  遵循与上述类似的原则。给定一个分区  $1$  进入  $k$  套  $1_1, \dots, 1_k$ ，我们定义  $k$  指标向量  $H = (H_{1j}, \dots, H_{nj})$  经过

$$H_{ij} = \begin{cases} 1/\sqrt{|1_j|} & \text{如果 } v_{-i} \in 1_j \\ 0 & \text{否则} \end{cases} \quad (1_j = 1, \dots, \tilde{n}_j = 1, \dots, k). \quad (5)$$

然后我们设置矩阵  $H \in \mathbb{R}^{n \times k}$  作为包含那些的矩阵  $k$  指标向量作为列。观察中的列  $H$  彼此正交，即  $H^T H = I$ 。在上一节的计算中我们可以看到

$$H_{-i}^T L H_{-i} = \frac{\text{切}(1_{-i}, 1_{-i})}{|1_{-i}|}.$$

此外，可以检查

$$H_{-i}^T L H_{-i} = (H^T L H)_{-i, -i}.$$

结合我们得到的这些事实

$$\text{比例切割}(1_1, \dots, 1_k) = \sum_{-i=1}^n H_{-i}^T L H_{-i} = \sum_{-i=1}^n (\text{高低})_{-i, -i} = \text{Tr}(H^T L H),$$

其中  $\text{Tr}$  表示矩阵的迹。所以最小化  $\text{RatioCut}(c_1, \dots, c_k)$  可以改写为

$$\min_{c_1, \dots, c_k} \text{Tr}(H^{-1} H) \text{ 受 } H^{-1} H = I, H \text{ 如方程式中所定义。} \quad (5).$$

与上面类似，我们现在通过允许矩阵的条目来放松问题  $H$  取任意实数值。那么松弛的问题就变成了：

$$\min_{H \in \mathbb{R}^{n \times k}} \text{Tr}(H^{-1} H) \text{ 受 } H^{-1} H = I.$$

这是迹最小化问题的标准形式，Rayleigh-Ritz 定理的一个版本（例如，参见 Lütkepohl, 1997 年的第 5.2.2.(6) 节）告诉我们解决方案是通过选择  $H$  作为包含第一个的矩阵  $k$  的特征向量作为列。我们可以看到矩阵  $H$  实际上是矩阵  $U$  在第 4 节中描述的非归一化谱聚类算法中使用。我们再次需要将实值解矩阵重新转换为离散分区。如上所述，标准方法是使用  $k$ -表示行上的算法  $U$ 。这导致了第 4 节中介绍的通用非归一化谱聚类算法。

### 5.3 逼近 Ncut

与用于 RatioCut 的技术非常相似的技术可用于导出归一化光谱聚类作为最小化 Ncut 的松弛。在这种情况下  $k=2$  我们定义聚类指标向量  $F$  经过

$$F_i = \begin{cases} \sqrt{\frac{\text{volume}(v)}{\text{volume}(V)}} & \text{如果 } v \in c_1 \\ -\sqrt{\frac{\text{volume}(v)}{\text{volume}(V)}} & \text{如果 } v \in c_2 \end{cases} \quad (6)$$

与上面类似，可以检查  $(Df)^T = 0, F^T Df = \text{volume}(V)$ , 和  $F^T F = \text{volume}(V)$  切割  $(c_1, c_2)$ . 因此我们可以用等价问题重写最小化 Ncut 的问题

$$\min_{c_1} F^T F \text{ 如果受制于 } F^T Df = \text{volume}(V), \text{ 自由度 } \perp 1 \text{ 个}, F^T Df = \text{volume}(V). \quad (7)$$

我们再次通过允许放松问题  $F$  取任意实数值：

$$\min_{F \in \mathbb{R}^n} F^T F \text{ 如果受制于 } Df^T \perp 1 \text{ 个}, F^T Df = \text{volume}(V). \quad (8)$$

现在我们代入  $G = \frac{1}{\sqrt{2}} F$ . 替换后，问题是

$$\min_{G \in \mathbb{R}^n} G^T G \text{ 受 } \frac{1}{\sqrt{2}} G^T Df = \frac{1}{\sqrt{2}} \text{volume}(V) \text{ 和 } \|G\|_2^2 = \text{volume}(V). \quad (9)$$

观察那个  $\frac{1}{\sqrt{2}} G^T Df = \frac{1}{\sqrt{2}} \text{volume}(V)$  是大号符号,  $\frac{1}{\sqrt{2}} \text{volume}(V)$  是第一个特征向量大号符号, 和  $\text{volume}(V)$  是一个常数。因此，问题 (9) 是标准 Rayleigh-Ritz 定理的形式，及其解  $G$  由第二个特征向量给出大号符号. 重新代入  $F = \sqrt{2} G$  并使用命题 3 我们看到  $F$  是第二个特征向量大号符号, 或等效的广义特征向量  $\lambda = \lambda_2$ .

对于发现的情况  $k > 2$  个集群，我们定义指标向量  $H = (H_{1j}, \dots, H_{nj})$  经过

$$H_{ij} = \begin{cases} \sqrt{\frac{\text{volume}(v)}{\text{volume}(V)}} & \text{如果 } v \in c_i \\ 0 & \text{否则} \end{cases} \quad (i=1, \dots, \tilde{n}_j=1, \dots, k). \quad (10)$$

图 2: Guattery 和 Miller (1998) 的蟑螂图。

然后我们设置矩阵  $H$  作为包含那些的矩阵  $k$  指标向量作为列。观察那个  $H$  不是我,  $H^T D H = I$ , 和  $H^T L H = \text{切}(\text{一个} \rightarrow \text{一个}, \text{一个} \rightarrow \text{一个}) / \text{体积}(\text{一个} \rightarrow \text{一个})$ . 所以我们可以写出最小化问题切割为

$$\min_{H \in \mathbb{R}^{n \times k}} \text{Tr}(H^T L H) \quad \text{受 } H^T D H = I, \quad H \text{ 如 (10)}.$$

放宽离散条件并代入  $H = \frac{1}{\sqrt{k}} U$  我们得到了松弛的问题

$$\min_{U \in \mathbb{R}^{n \times k}} \text{Tr}(U^T L U) \quad \text{受 } U^T D U = I. \quad (11)$$

同样, 这是由矩阵解决的标准迹线最小化问题  $U$  其中包含第一个  $k$  的特征向量  $D$  符号作为列。重新代入  $H = \frac{1}{\sqrt{k}} U$  并使用命题 3 我们看到解决方案  $H$  由第一个组成  $k$  矩阵的特征向量  $D$  写字, 或第一个  $k$  的广义特征向量  $U = \lambda U$ . 这产生了根据 Shi 和 Malik (2000) 的归一化谱聚类算法。

## 5.4 对放松方法的评论

关于谱聚类的这种推导, 我们应该做出几点评论。最重要的是, 与精确解相比, 松弛问题的解的质量没有任何保证。也就是说, 如果一个  $1 \rightarrow \dots$ , 一个  $k$  是最小化 RatioCut 的精确解, 并且  $Z_1 \rightarrow \dots, Z_k$  是非归一化谱聚类构造的解, 则  $\text{RatioCut}(Z_1 \rightarrow \dots, Z_k) - \text{比例切割}(\text{一个} \rightarrow \text{一个}, \text{一个} \rightarrow \text{一个})$  可以任意大。在 Guattery 和 Miller (1998) 中可以找到这方面的几个例子。例如, 作者考虑了一类非常简单的图, 称为“蟑螂图”。这些图基本上看起来像一个梯子, 去掉了一些边缘, 见图 2。显然, 理想的  $\text{RatioCut } k=2$  只是通过垂直切割来切割梯子, 这样  $\text{一个} = \{v_1 \rightarrow \dots, v_k, v_{2 \rightarrow k+1}, \dots, v_{3 \rightarrow k}\}$  和  $\text{一个} = \{v_{k+1}, \dots, v_{2 \rightarrow k}, v_{3 \rightarrow k+1}, \dots, v_{4 \rightarrow k}\}$ . 这种剪裁与  $\text{一个} \rightarrow \text{一个} \rightarrow 2 \rightarrow k$  和  $\text{一个} \rightarrow \text{一个} = 2$ . 然而, 通过研究蟑螂图的非归一化图拉普拉斯算子的第二个特征向量的性质, 作者证明了非归一化谱聚类总是通过阶梯水平切割, 构造集合  $Z = \{v_1 \rightarrow \dots, v_{2 \rightarrow k}\}$  和  $Z = \{v_{2 \rightarrow k+1}, \dots, v_{4 \rightarrow k}\}$ . 这也会导致平衡切割, 但现在我们切割  $k$  边而不仅仅是 2。所以  $\text{RatioCut}(\text{一个}, \text{一个}) = 2/k$ , 而  $\text{RatioCut}(Z, Z) = 1$ 。这意味着相对于最优切割, 谱聚类得到的 RatioCut 值为  $k/2$  倍, 这是顺序中的一个因素  $n$ . 其他几篇论文研究了由谱聚类构建的聚类的质量, 例如 Spielman 和 Teng (1996) (针对非归一化谱聚类) 和 Kannan、Vempala 和 Vetta (2004) (针对归一化谱聚类)。通常, 已知不存在将平衡图切割近似为常数因子的有效算法。相反, 这个近似问题本身可能是 NP 问题 (Bui 和 Jones, 1992)。

当然，我们上面讨论的松弛并不是唯一的。例如，Bie 和 Cristianini (2006) 推导了导致半定规划的完全不同的松弛，可能还有许多其他有用的松弛。谱松弛之所以如此吸引人，并不是因为它能带来特别好的解决方案。它的流行主要是因为它导致了一个标准的线性代数问题，它很容易解决。

## 6 随机游走的观点

解释谱聚类的另一个论据是基于相似度图上的随机游走。图上的随机游走是一个从一个顶点随机跳到另一个顶点的随机过程。我们将在下面看到，谱聚类可以解释为试图找到图的分区，使得随机游走在同一簇内停留很长时间，并且很少在簇之间跳跃。直觉上这是有道理的，特别是结合上一节的图切割解释：具有低切割的平衡分区也将具有随机游走在集群之间跳跃的机会不多的特性。对于一般随机游走的背景阅读，我们参考 Norris (1997) 和 Brémaud (1999)，对于图上的随机游走，我们推荐 Aldous 和 Fill (准备中) 和 Lovász (1993)。正式地， $d$ -步到顶点  $v_j$  与边权重成正比  $w_{ij}$  并由  $p_{ij} = w_{ij}/d$ 。转移矩阵  $P = (p_{ij})_{i,j=1,\dots,n}$  因此，随机游走的定义为

$$P = D^{-1}W.$$

如果图是连通的且非二分的，则随机游走总是具有唯一的平稳分布  $\pi = (\pi_1, \dots, \pi_n)$ ，在哪里  $\pi_i = d_i / \text{体积}(V)$ 。显然两者之间存在紧密的关系。大写字母  $P$  作为大写字母  $\pi$  的结果， $\lambda$  是特征值大写字母与特征向量  $\phi$  当且仅当  $1 - \lambda$  是特征值  $P$  与特征向量  $\phi$ 。众所周知，图的许多属性可以用相应的随机游走转移矩阵来表示  $P$ ，参见 Lovász (1993) 的概述。从这个角度来看，最大的特征向量并不奇怪  $P$  和最小的特征向量大写字母可用于描述图的簇属性。

### 随机游走和 Ncut

Meila 和 Shi (2001) 观察到 Ncut 和随机游走的转移概率之间的形式等价性。

提案 5 (切割通过转移概率) 让  $G$  连接且非双向。假设我们运行随机游走  $(X_t)_{t \in \mathbb{N}}$  从  $x_0$  在平稳分布  $\pi$  中。不相交的子集  $A, B \subset V$ ，表示为  $P(B|A) := P(X_1 \in B | X_0 \in A)$ 。然后：作为

$$\text{切割}(A, B) = P(A|B) + P(B|A).$$

证明。首先观察到

$$\begin{aligned} P(X_0 \in A, X_1 \in B) &= \sum_{i \in A, j \in B} P(X_0 = i, X_1 = j) = \sum_{i \in A, j \in B} \pi_i p_{ij} \\ &= \sum_{i \in A, j \in B} \frac{d_i w_{ij}}{\text{体积}(V) d_i} = \frac{1}{\text{体积}(V)} \sum_{i \in A, j \in B} w_{ij}. \end{aligned}$$

使用这个我们得到

$$P(X_1 \in B | X_0 \in \alpha) = \frac{P(X_0 \in \alpha, X_1 \in B)}{P(X_0 \in \alpha)} = \frac{\sum_{\beta \in Z} P(X_0 \in \alpha, X_1 \in \beta)}{\sum_{\beta \in Z} P(X_0 \in \alpha, X_1 \in \beta)} = \frac{\sum_{\beta \in Z} W_{\alpha\beta}}{\sum_{\beta \in Z} W_{\alpha\beta}}.$$

现在命题直接遵循 Ncut 的定义。

2个

这个命题导致对 Ncut 的一个很好的解释，因此对归一化谱聚类。它告诉我们，当最小化 Ncut 时，我们实际上是在图形中寻找切割，这样随机游走很少从一个至一个反之亦然。

## 通勤距离

随机游走和图拉普拉斯算子之间的第二个联系可以通过图上的通勤距离来建立。通勤距离（也称为阻力距离） $C_{ij}$ 两个顶点之间  $v_i$  和  $v_j$  是随机游走从顶点  $v_i$  行进所需的预期时间  $v_i$  到顶点  $v_j$  和返回 (Lovász, 1993 年; Aldous 和 Fill, 准备中)。通勤距离有几个很好的特性，这使得它对机器学习特别有吸引力。与图上的最短路径距离相反，如果有许多不同的从顶点到达的短路径，则两个顶点之间的通勤距离会减少  $v_i$  到顶点  $v_j$ 。因此，通勤距离不是仅仅寻找一条最短路径，而是着眼于的一组最短路径。在图中由短路径连接且位于图中相同高密度区域的点被认为比由短路径连接但位于图中不同高密度区域的点彼此更接近。从这个意义上说，通勤距离似乎特别适用于聚类目的。

值得注意的是，可以借助广义逆（也称为伪逆或 Moore-Penrose 逆）计算图上的通勤距离。图拉普拉斯算子  $L$  下面我们记  $e_i = (0, \dots, 0, 1, 0, \dots, 0)$  作为  $i$ -th 单位向量。定义广义逆  $L^+$ ，回想一下命题 1 的矩阵  $L$  可以分解为  $L = U\Lambda U^T$  在哪里  $U$  是包含所有特征向量作为列的矩阵， $\Lambda$  是具有特征值的对角矩阵  $\lambda_1, \dots, \lambda_n$  在对角线上。由于至少一个特征值是 0，矩阵  $L$  不可逆。相反，我们定义它的广义

逆为  $L^+ := U\Lambda^+ U^T$  其中矩阵  $\Lambda^+$  是对角元素为 1 的对角矩阵  $\lambda_i^{-1}$  如果  $\lambda_i \neq 0$  和 0 如果  $\lambda_i = 0$ 。因此  $L^+$  可以计算为  $L^+ = \sum_{i=1}^n \frac{1}{\lambda_i} u_i u_i^T$ 。矩阵  $L^+$  是积极的半定和对称。对于进一步的属性  $L^+$  参见 Gutman 和 Xiao (2004)。

提案 6（通勤距离）让  $G=(V,E)$  连通的无向图。表示为  $C_{ij}$

顶点之间的通勤距离  $v_i$  和顶点  $v_j$ ，通过  $L^+$  的  $L$ 。然后我们有：

$i, j = 1, \dots, n$  广义逆

$$C_{ij} = \text{卷}(V)(L^+ - 2\text{个升}_{ij} + \text{升}_{ii}) = \text{卷}(V)(e_i - e_j)(L^+)(e_i - e_j).$$

这一结果已由 Klein 和 Randic (1993) 发表，其中已通过电气网络理论的方法得到证明。有关使用第一步分析进行随机游走的证明，请参阅 Fouss、Pirrotte、Renders 和 Saerens (2007)。在图拉普拉斯算子的帮助下，还存在其他表达通勤距离的方法。例如，根据归一化拉普拉斯算子的特征向量的方法  $L^+$  符号

可以在 Lovász (1993) 中作为推论 3.2 找到，以及一种借助于某些子矩阵的行列式计算通勤距离的方法  $L^+$  可以在 Bapat、Gutman 和 Xiao (2003) 中找到。

提案 6 具有重要的意义。这表明  $C_{ij}$  可以被认为是欧几里得图的顶点上的距离函数。这意味着我们可以构建一个嵌入



映射顶点  $v$  到点图的  $z \in \mathbb{R}^n$  这样点之间的欧几里得距离  $\|z - z'\|_2$  与图表上的通勤距离一致。这工作如下。作为矩阵  $A$

是半正定和对称的，它在  $\mathbb{R}^n$  (或者更正式地说，它在子空间上导出一个内积  $\mathbb{R}^n$  垂直于矢量  $\mathbf{1}$ )。现在选择  $z$  作为点  $R$  对应于  $A$  矩阵的第行  $\tilde{u}(\Lambda)_{1 \times 2 \times n}$ 。那么，根据命题 6

并通过建造  $A$  我们有那个  $\langle z - z', z_j \rangle = \text{电子} \cdot \text{—世大号} + \text{电子和 } C_j = \text{体积} \cdot (V) / \|z - z'\|_2$ 。

非归一化谱聚类中使用的嵌入与通勤时间嵌入有关，但不完全相同。在谱聚类中，我们将图的顶点映射到行上是  $A$  矩阵的  $\tilde{u}$ ，而通勤时间嵌入映射行上的顶点  $z$  到  $A$  矩阵的  $(\Lambda)_{1 \times 2 \times n} \tilde{u}$ 。也就是说，与条目相比是  $A$ ，条目  $z$  由的反特征值另外缩放大号。此外，在谱聚类中，我们只取第一个  $k$  矩阵的列，而通勤时间嵌入采用所有列。几位作者现在试图证明为什么是  $A$  和  $z$  毕竟并没有太大的不同，并且有点放弃，事实是谱聚类基于欧几里德距离构建聚类是  $A$  可以解释为基于通勤距离构建图中顶点的集群。但是，请注意，这两种方法可能有很大不同。例如，在图形包含的最佳情况下  $k$  断开连接的组件，第一个  $k$  的特征值  $A$  根据命题 2 为 0，第一个  $k$  列的  $\tilde{u}$  由聚类指标向量组成。然而，第一  $k$  矩阵的列  $(\Lambda)_{1 \times 2 \times n} \tilde{u}$  仅由零组成，作为第一个  $k$  的对角元素  $t$  为 0。在这种情况下，第一个中包含的信息  $k$  列的  $\tilde{u}$  在矩阵中完全被忽略  $(\Lambda)_{1 \times 2 \times n} \tilde{u}$ ，以及矩阵的所有非零元素  $(\Lambda)_{1 \times 2 \times n} \tilde{u}$  可以在列中找到  $k+1$  到  $n$  在光谱聚类中没有考虑，它会丢弃所有这些列。另一方面，如果底层图是连通的，则不会出现这些问题。在这种情况下，唯一具有特征值 0 的特征向量是常数 1 向量，在这两种情况下都可以忽略它。对应于小特征值的特征向量  $\lambda$

的大号然后在矩阵中强调  $(\Lambda)_{1 \times 2 \times n} \tilde{u}$  因为它们乘以  $\lambda$   $—世=1 \times \lambda - 世$ 。在这种情况下——通勤时间嵌入和频谱嵌入做类似的事情可能是真的。

总而言之，通勤时间距离似乎是一个有用的直觉，但如果不做进一步的假设，谱聚类与通勤距离之间的关系就相当松散。这些关系可能会收紧，例如，如果相似函数是严格正定的。然而，我们还没有看到关于这一点的精确数学陈述。

## 7 微扰理论观点

微扰理论研究矩阵的特征值和特征向量如何变化的问题一个如果我们添加一个小的扰动就会改变  $H$ ，即我们考虑扰动矩阵  $A = A + H$ 。大多数微扰定理指出特征值或特征向量之间的一定距离一个和一个以常数乘以范数为界  $H$ 。该常数通常取决于我们正在查看的特征值，以及该特征值与光谱的其余部分分开的距离（有关正式声明，请参见下文）。谱聚类的理由如下：让我们首先考虑“理想情况”，其中簇间相似性恰好为 0。我们在第 3 节中看到，第一个  $k$  的特征向量  $A$  或者  $A$  写字是聚类的指标向量。在这种情况下，点是  $z \in \mathbb{R}^k$

在谱聚类算法中构建的形式为  $(0, \dots, 0, 1, \dots, 0)$  其中 1 的位置表示该点所属的连通分量。特别是，所有是  $A$  属于同一个连通分量重合。这  $k$ -means 算法将通过在每个点上放置一个中心点  $(0, \dots, 0, 1, \dots, 0) \cdot \in \mathbb{R}^k$ 。在“近乎理想的情况”中，我们仍然有不同的集群，但集群之间的相似性不完全为 0，我们认为拉普拉斯矩阵是理想情况下的扰动版本。然后微扰理论告诉我们，特征向量将非常接近理想的指标向量。积分是  $A$  也许不会



完全符合  $(0, \dots, 0, 1, 0, \dots, 0)$ ；但这样做会导致一些小的误差项。因此，如果扰动不是太大，则  $k$ -means 算法仍然会将组彼此分开。

## 7.1 形式微扰论证

谱聚类扰动方法的正式基础是矩阵扰动理论中的戴维斯-卡汉定理。该定理限制了扰动下对称矩阵的特征空间之间的差异。我们陈述这些结果是为了完整性，但对于背景阅读，我们参考了 Stewart 和 Sun (1990) 的第 V 节和 Bhatia (1997) 的第 VII.3 节。在微扰理论中，子空间之间的距离通常使用“标准角”（也称为“主角”）来测量。要定义主角，让  $V_1$  和  $V_2$  是两个  $p$ -维子空间  $\mathbb{R}^d$ ，和  $V_1$  和  $V_2$  两个矩阵使得它们的列形成标准正交系统  $V_1$  和  $V_2$ ，分别。然后

余弦  $\cos \Theta$ —主角  $\Theta$ —是奇异值  $V_1^T V_2$ 。为了  $p=1$ ，如此定义规范角与角的法线定义一致。也可以定义标准角，如果  $V_1$  和  $V_2$  不具有相同的维度，请参见 Stewart 和 Sun (1990) 的第 V 节、Bhatia (1997) 的第 VII.3 节或 Golub 和 Van Loan (1996) 的第 12.4.3 节。矩阵  $\sin \Theta(V_1, V_2)$  将表示对角矩阵，其正弦值位于对角线上。

定理 7 (戴维斯-卡汉) 让  $A \in \mathbb{R}^{n \times n}$  是对称矩阵，让  $\|\cdot\|$  分别是 Frobenius 范数或矩阵的双范数。考虑  $\tilde{A} := A + H$  作为一个扰动的版本  $A$ 。让  $\sigma_1 \in \mathbb{C}$  成为一个区间。用  $\sigma$  表示  $\sigma_1$  的一个特征值集  $\sigma_1$  包含在  $\sigma_1$  中，通过  $V_1$  对应于所有这些特征值的特征空间（更正式地说， $V_1$  是由  $\sigma$  引起的光谱投影的图像  $V_1(\sigma_1)$  一个）。用  $\sigma$  表示  $\sigma_1$  的一个和  $\tilde{V}_1$  的类似数量。定义之间的距离  $d(V_1, \tilde{V}_1)$  和光谱  $\sigma_1$  在外面  $\sigma_1$  作为

$$\delta = \min_{\lambda \in \sigma_1} \min_{\lambda' \in \sigma_1'} |\lambda - \lambda'|, \quad \lambda \in \sigma_1, \lambda' \in \sigma_1'.$$

那么距离  $d(V_1, \tilde{V}_1) := \|\sin \Theta(V_1, \tilde{V}_1)\|$  两个子空间之间  $V_1$  和  $\tilde{V}_1$  受限于

$$d(V_1, \tilde{V}_1) \leq \frac{\|H\|}{\delta}.$$

有关该定理的讨论和证明，请参见 Stewart 和 Sun (1990) 的第 V.3 节。让我们尝试解密这个定理，为了简单起见，在非标准化拉普拉斯算子的情况下（对于标准化拉普拉斯算子，它的工作方式类似）。矩阵  $A$  一个将对应于拉普拉斯图  $G$  在图表具有的理想情况下  $k$  连接的组件。矩阵  $A$  一个对应于一个扰动的情况，其中由于噪声  $G$  图中的组件不再完全断开，而是仅由少数权重较低的边连接。我们将这种情况的相应图拉普拉斯算子表示为  $\tilde{A}$ 。对于谱聚类，我们需要首先考虑  $k$  的特征值和特征向量  $V_k$  表示的特征值  $\lambda_1, \dots, \lambda_k$  和不安的拉普拉斯算子

$\tilde{A}$  经过  $\lambda_1, \dots, \lambda_k$ 。选择间隔  $\sigma_1$  现在是关键点。我们想这样选择它都是第一个  $k$  的特征值  $\lambda_k$  和第一个  $k$  的特征值  $\lambda_{k+1}$  包含在  $\sigma_1$  中。扰动越小越容易  $H = \tilde{A} - A$  本征隙越大  $|\lambda_k - \lambda_{k+1}|$  是。如果我们设法找到这样一个集合，那么 Davis-Kahan 定理告诉我们对应于第一个的特征空间  $k$  理想矩阵的特征值  $\lambda_k$  和第一个  $k$  扰动矩阵的特征值  $\tilde{\lambda}_k$  彼此非常接近，也就是说它们的距离以  $\|H\|/\delta$ 。然后，由于理想情况下的特征向量在连通分量上是分段常数，因此在扰动情况下也大致如此。“大约”有多好取决于扰动的范数  $\|H\|$  和距离  $\delta$  之间  $\sigma_1$  和  $(k+1)$  的特征向量  $\lambda_{k+1}$ 。如果设置  $\sigma_1$  已被选为区间  $[\lambda_k, \lambda_{k+1}]$ ，然后  $\delta$  与光谱间隙重合  $|\lambda_{k+1} - \lambda_k|$ 。从定理可以看出，这个特征间隙越大，理想情况和扰动情况的特征向量越接近，谱聚类效果越好。下面我们将看到本征间隙的大小也可以用在

不同的上下文作为谱聚类的质量标准，即在选择数字时 $k$ 要构建的集群。

如果扰动 $H$ 太大或特征间隙太小，我们可能找不到集合 $\mathcal{V}_1$ 这样无论是第一个 $k$ 的特征值 $\lambda_1$ 和 $\mathcal{V}_1$ 包含在 $\mathcal{V}_1$ 。在这种情况下，我们需要通过选择集合来做出折衷 $\mathcal{V}_1$ 包含第一个 $k$ 的特征值 $\lambda_1$ ，但可能有一些或多或少的特征值 $\lambda_i$ 。然后定理的陈述变得更弱，因为我们要么不比较对应于第一个的特征空间 $k$ 的特征向量 $\mathbf{v}_1$ 和 $\mathcal{V}_1$ ，但与第一个对应的特征空间 $k$ 的特征向量 $\mathbf{v}_1$ 和第一个 $k$ 的特征向量 $\mathbf{v}_1$ （在哪里 $k$ 是特征值的数量 $L$ 包含在 $\mathcal{V}_1$ ）。或者，它可能发生 $\delta$ 变得如此之小，以至于两者之间距离的界限 $d(\mathcal{V}_1, \tilde{\mathcal{V}}_1)$ 爆炸太多以至于它变得毫无用处。

## 7.2 关于微扰方法的评论

在使用微扰理论来证明基于矩阵特征向量的聚类算法时需要谨慎一些。一般来说，任何分块对角对称矩阵具有这样的特性，即存在特征向量的基，这些特征向量在各个块外为零，在块内为实值。例如，基于这个论点，几位作者使用了相似矩阵的特征向量 $\mathbf{v}_1$ 或邻接矩阵 $W$ 发现集群。然而，在完全分离的簇的理想情况下是块对角线可以被认为是成功使用特征向量的必要条件，但不是充分条件。至少还应满足两个属性：

首先，我们需要确保 $\mathbf{v}_1$ 的特征值和特征向量是有意义的。在拉普拉斯算子的情况下，这总是正确的，因为我们知道任何连通分量都恰好拥有一个特征值为0的特征向量。因此，如果图有 $k$ 连接的组件，我们采取第一个 $k$ 拉普拉斯算子的特征向量，那么我们就知道每个分量只有一个特征向量。但是，对于其他矩阵可能并非如此，例如 $\mathbf{v}_1$ 或者 $W$ 。例如，可能是块对角相似矩阵的两个最大特征值 $\lambda_1$ 来自同一个街区。在这种情况下，如果我们采取第一个 $k$ 的特征向量 $\mathbf{v}_1$ ，一些块将出现多次，而另一些块我们将完全错过（除非我们采取某些预防措施）。这就是为什么使用特征向量的原因 $\mathbf{v}_1$ 或者 $W$ 不鼓励聚类。

第二个属性是，在理想情况下，组件上的特征向量的条目应该“安全地远离”0。假设第一个连接的组件上的特征向量有一个条目 $v_{1,i} > 0$ 在位置 $i$ 。在理想情况下，此条目非零这一事实表明相应的点 $i$ 属于第一个集群。反过来，如果一个点 $j$ 不属于集群1，那么在理想情况下应该是这样的 $v_{1,j} = 0$ 。现在考虑相同的情况，但数据受到扰动。扰动的特征向量 $\tilde{\mathbf{v}}_1$ 通常不再有任何非零分量；但如果噪声不是太大，那么微扰理论告诉我们条目 $\tilde{v}_{1,i}$ 和 $\tilde{v}_{1,j}$ 仍然“接近”它们的原始值 $v_{1,i}$ 和 $v_{1,j}$ 。所以两个条目 $\tilde{v}_{1,i}$ 和 $\tilde{v}_{1,j}$ 会取一些小的值，比如 $\epsilon_1$ 和 $\epsilon_2$ 。实际上，如果这些值非常小，我们就不清楚我们应该如何解释这种情况。要么我们相信小条目 $\tilde{v}$ 表示这些点不属于第一个集群（然后错误分类第一个数据点 $i$ ），或者我们认为这些条目已经表明了类成员并将两个点都分类到第一个集群（它错误地分类了点 $j$ ）。

对于两个矩阵 $\mathbf{v}_1$ 和 $\mathbf{v}_2$ ，理想情况下的特征向量是指示向量，所以不会出现上述第二个问题。然而，这对于矩阵来说是不正确的 $\mathbf{v}_1$ 符号，用于Ng等人的归一化谱聚类算法。(2002)。即使在理想情况下，该矩阵的特征向量也为 $\mathbf{v}_1/\sqrt{1}$ 。如果顶点的度数相差很大，特别是如果存在度数很低的顶点，则特征向量中的对应项非常小。为了解决上述问题，Ng等人的算法中的行归一化步骤。(2002)开始发挥作用。在理想情况下，矩阵 $\tilde{L}$ 在算法中只有一个

每行非零条目。行归一化后，矩阵 $U$ 在 Ng 等人的算法中。(2002) 然后由聚类指标向量组成。但是请注意，这在实践中可能并不总是正确。假设我们有  $\tilde{U}_{-i}$ ,  $1 \leq i \leq n$  和  $\sqrt{\frac{1}{1 + \epsilon_2 \epsilon_1}}$ 。

如果我们现在将  $\tilde{U}_{-i}$  乘以  $\frac{1}{1 + \epsilon_2 \epsilon_1}$  并且变得相当大。我们现在遇到与上述类似的问题：两个点很可能被归类到同一个集群中，即使它们属于不同的集群。这个论点表明使用矩阵的谱聚类大符号如果特征向量包含特别小的条目，则可能会出现这个问题。另一方面，请注意，特征向量中的此类小条目仅在某些顶点的度数特别低时才会出现（如特征向量大符号由  $\sqrt{1/\epsilon_2}$  个一个）。有人可能会争辩说，在这种情况下，数据点无论如何都应该被视为异常值，那么该点最终会出现在哪个集群中并不重要。

总而言之，结论是非归一化谱聚类和归一化谱聚类都具有大符号字微扰理论方法很好地证明了这一点。归一化谱聚类大符号也可以用微扰理论来证明，但是如果图中包含度数非常低的顶点，则应该更加小心。

## 8 实用细节

在本节中，我们将简要讨论实际实施谱聚类时出现的一些问题。有几个选择和要设置的参数。但是，本节中的讨论主要是为了提高对发生的一般问题的认识。为了深入研究各种现实世界任务的谱聚类行为，我们参考了文献。

### 8.1 构建相似度图

为光谱聚类构建相似度图并不是一项微不足道的任务，并且对各种结构的理论含义知之甚少。

#### 相似函数本身

在我们甚至可以考虑构建相似度图之前，我们需要在数据上定义一个相似度函数。因为我们稍后要构建一个邻域图，所以我们需要确保由这个相似函数导出的局部邻域是“有意义的”。这意味着我们需要确保被相似度函数认为“非常相似”的点在数据来源的应用程序中也密切相关。例如，在构建文本文档之间的相似度函数时，检查具有高相似度得分的文档是否确实属于同一文本类别是有意义的。相似性函数的全局“长程”行为对于谱聚类并不那么重要——两个数据点的相似性得分是 0.01 还是 0.001 并不重要，比如说，因为无论如何我们都不会在相似图中连接这两个点。在数据点存在于欧几里德空间中的常见情况下  $R^d$ ，一个合理的默认候选者是高斯相似函数  $s(x_i, x_j) = \exp(-\|x_i - x_j\|_2^2 / (2\sigma^2))$  (当然我们需要选择参数  $\sigma$  在这里，见下文)。归根结底，相似度函数的选择取决于数据来自的领域，无法给出一般性建议。

#### 哪种类型的相似度图

必须做出的下一个选择涉及要使用的图形类型，例如  $k$ -最近的邻居或  $\epsilon$ -邻域图。让我们使用图 3 中显示的玩具示例来说明不同图形的行为。作为基础分布，我们选择分布  $R_2$  个和

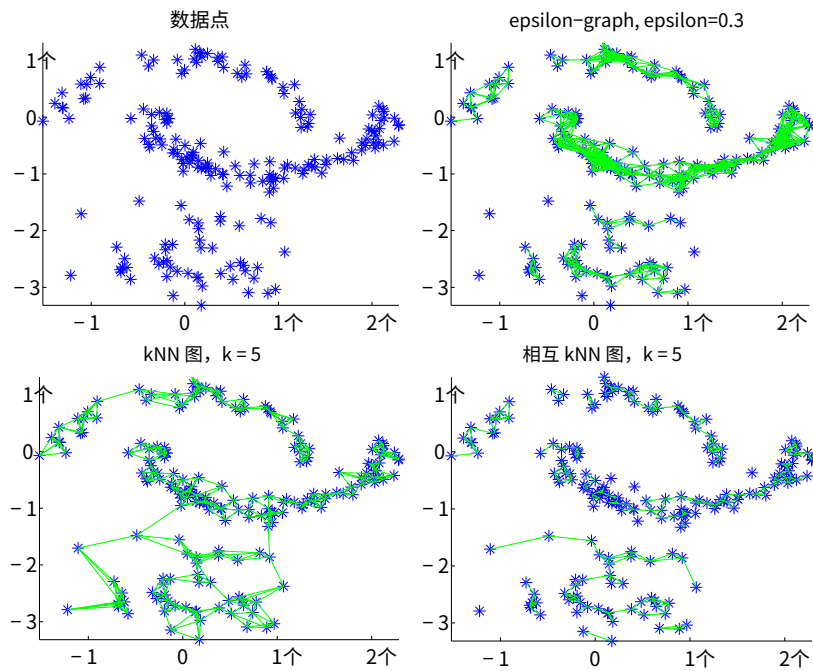


图 3：不同的相似度图，详见正文。

三个星团：两个“月亮”和一个高斯星团。底部卫星的密度被选择为大于顶部卫星的密度。图 3 中的左上面板显示了从该分布中提取的样本。接下来的三个面板显示了此样本的不同相似度图。

在里面  $\epsilon$ -neighborhood graph，我们可以看到很难选择一个有用的参数  $\epsilon$ 。和  $\epsilon=0.3$  如图，中间moon上的点已经很紧密的连在一起了，而Gaussian上的点几乎没有连在一起。如果我们有“不同尺度”的数据，那么这个问题总是会发生，即数据点之间的距离在空间的不同区域是不同的。

这  $k$ - 另一方面，最近邻图可以连接“不同尺度”的点。我们可以看到低密度高斯中的点与高密度月球中的点相连。这是一般属性  $k$ -最近的邻居图，这可能非常有用。我们还可以看到  $k$ - 如果存在彼此相距相当远的高密度区域，则最近邻图可能会分成几个不相连的组件。本例中的两个卫星就是这种情况。

相互的  $k$ - 最近邻图具有这样的特性，即它倾向于连接恒定密度区域内的点，但不连接不同密度的区域。所以相互的  $k$ - 最近邻图可以被认为是“介于”  $\epsilon$ -邻域图和  $k$ - 最近邻图。它能够在不同的尺度上起作用，但不会将这些尺度相互混合。因此，相互  $k$ - 如果我们想要检测不同密度的集群，最近邻图似乎特别适合。

全连接图经常与高斯相似函数结合使用  $\text{sim}(X_i, X_j) = \exp(-\|X_i - X_j\|_2^2 / (2\sigma^2))$ 。这里的参数  $\sigma$  起着与参数相似的作用  $\epsilon$  在里面  $\epsilon$ -邻域图。局部邻域中的点与相对较高的权重相连，而  $r$  个远离点之间的边具有正的但可以忽略不计的权重。然而，由此产生的

一个

相似矩阵不是稀疏矩阵。

作为一般性建议，我们建议与 $k$ -最近邻图作为首选。使用起来很简单，导致稀疏邻接矩阵 $W$ ，并且根据我们的经验，与其他图表相比，它更不容易受到参数选择不当的影响。

#### 相似度图的参数

一旦决定了相似图的类型，就必须选择其连接参数 $k$ 或者 $\epsilon$ ，分别。不幸的是，几乎没有任何理论结果可以指导我们完成这项任务。一般来说，如果相似度图包含的连通分量多于我们要求算法检测的聚类数量，那么谱聚类将平凡地将连通分量作为聚类返回。除非完全确定那些连接的组件是正确的集群，否则应该确保相似图是连接的，或者只包含“很少”的连接组件和很少或没有孤立的顶点。关于如何实现随机图的连通性有很多理论结果，但所有这些结果都只在样本量的限制下成立 $n \rightarrow \infty$ 。例如，众所周知，对于 $n$ 从一些底层密度中提取的数据点具有连接支持 $R_d$ ，这 $k$ -最近邻图和相互 $k$ -如果我们选择，最近邻图将被连接 $k$ 按照日志的顺序 $(n)$ （例如，Brito、Chavez、Quiroz和Yukich，1997）。类似的论点表明参数 $\epsilon$ 在里面 $\epsilon$ -邻域图必须选择为 $(\log(n)/n)^d$ 以保证在极限范围内的连通性(Penrose, 1999)。尽管具有理论上的意义，但所有这些结果并不能真正帮助我们进行选择 $k$ 在有限样本上。

现在让我们给出一些经验法则。使用时 $k$ -最近邻图，那么应该选择连通性参数，使得生成的图是连通的，或者至少具有比我们想要检测的簇少得多的连通分量。对于小型或中型图形，可以“步行”进行尝试。对于非常大的图，第一个近似值可能是选择 $k$ 按照日志的顺序 $(n)$ ，正如渐近连通性结果所建议的那样。

对于相互 $k$ -最近邻图，我们不得不承认我们有点迷失了经验法则。互惠互利的优势 $k$ -与标准图相比的最近邻图 $k$ -最近邻图是它倾向于不连接不同密度的区域。虽然如果存在由单独的高密度区域引起的清晰聚类，这可能会很好，但在不太明显的情况下这可能会造成伤害，因为图中不连贯的部分将始终被谱聚类选择为聚类。一般而言，可以观察到相互 $k$ -最近邻图的边比标准图少得多 $k$ -相同参数的最近邻图 $k$ 。这建议选择 $k$ 对共同体来说要大得多 $k$ -最近的邻居图比一个人会做的标准 $k$ -最近邻图。但是，要利用相互的属性 $k$ -最近邻图不连接不同密度的区域，因此有必要允许图中有几个“有意义的”断开连接的部分。不幸的是，我们不知道有任何通用的启发式方法来选择参数 $k$ 这样就可以实现。

为了 $\epsilon$ -neighborhood graph，建议选择 $\epsilon$ 这样生成的图就安全地连接起来了。确定最小值 $\epsilon$ 图连接的地方很简单：必须选择 $\epsilon$ 作为数据点上全连接图的最小生成树中最长边的长度。后者可以通过任何最小生成树算法轻松确定。但是，请注意，当数据包含异常值时，该启发式方法将选择 $\epsilon$ 如此之大以至于即使异常值也与其他数据相关联。当数据包含几个彼此相距很远的紧密集群时，会发生类似的效果。在这两种情况下， $\epsilon$ 将选择太大以反映数据最重要部分的规模。

最后，如果将全连接图与可缩放的相似度函数一起使用

本身，例如高斯相似函数，那么应该选择相似函数的尺度，使得结果图具有与相应的相似属性  $k$ -最近的邻居或  $\varepsilon$ -邻域图会有。需要确保对于大多数数据点，相似度显著大于 0 的邻居集“不太小也不太大”。特别是，对于高斯相似度函数，经常使用几个经验法则。例如，可以选择  $\sigma$  按照一个点到它的平均距离的顺序  $k$ -th 最近的邻居，其中  $k$  与上面类似地选择（例如， $k \sim \log(n) + 1$ ）。另一种方法是确定  $\varepsilon$  通过上面描述的最小生成树启发式，然后选择  $\sigma = \varepsilon$ 。但请注意，所有这些经验法则都是非常临时的，并且取决于手头的给定数据及其点间距离的分布，它们可能根本不起作用。

一般来说，经验表明谱聚类对相似图的变化及其参数的选择非常敏感。不幸的是，据我们所知，还没有系统的研究来调查相似性图及其参数对聚类的影响，并提出合理的经验法则。上述建议都没有坚实的理论基础。寻找具有理论依据的规则应被视为未来研究的一个有趣且重要的课题。

## 8.2 计算特征向量

要在实践中实现谱聚类，必须计算第一个  $k$  潜在大图拉普拉斯矩阵的特征向量。幸运的是，如果我们使用  $k$ -最近邻图或  $\varepsilon$ -邻域图，那么所有这些矩阵都是稀疏的。存在计算稀疏矩阵的第一个特征向量的有效方法，最流行的方法是幂方法或 Krylov 子空间方法，例如 Lanczos 方法（Golub 和 Van Loan，1996）。这些算法的收敛速度取决于特征间隙（也称为谱间隙）的大小  $|\lambda_k - \lambda_{k+1}|$ 。这个 eigengap 越大，算法计算第一个的速度越快  $k$  特征向量收敛。

请注意，如果所考虑的特征值之一具有大于一的多重性，则会出现一般问题。例如，在理想情况下  $k$  断开连接的簇，特征值 0 具有多重性  $k$ 。正如我们所见，在这种情况下，特征空间由  $k$  聚类指标向量。但不幸的是，由数值特征求解器计算的向量不一定会收敛到那些特定的向量。相反，它们只是收敛到本征空间的某个标准正交基，并且通常取决于算法收敛到哪个基的实现细节。但是这个

毕竟还不错。请注意， $\sum_k$  集群指示符向量跨越的空间中的所有向量 1 个 一个一世  
有形式  $\sum_k \alpha_k \mathbf{1}_k$  对于一些系数  $\alpha_k$  一个一世，也就是说，它们在  
集群。因此，特征求解器返回的向量仍然编码了有关聚类的信息，然后可以由  $k$  意味着重建集群的算法。

## 8.3 簇数

选择号码  $k$  聚类问题是所有聚类算法的一个普遍问题，针对这个问题已经设计出各种或多或少成功的方法。在基于模型的聚类设置中，存在从数据中选择聚类数量的合理标准。这些标准通常基于数据的对数似然性，然后可以以频率论或贝叶斯方式处理，例如参见 Fraley 和 Raftery (2002)。在没有或很少对基础模型做出假设的情况下，可以使用大量不同的指标来选择聚类的数量。示例包括临时测量，例如集群内和集群间相似性的比率、信息论标准（Still 和 Bialek，2004 年）、差距统计（Tibshirani、Walther 和 Hastie，2001 年），以及稳定性方法（宾虚，Elisseeff 和 Guyon，2002 年；兰格，罗斯，

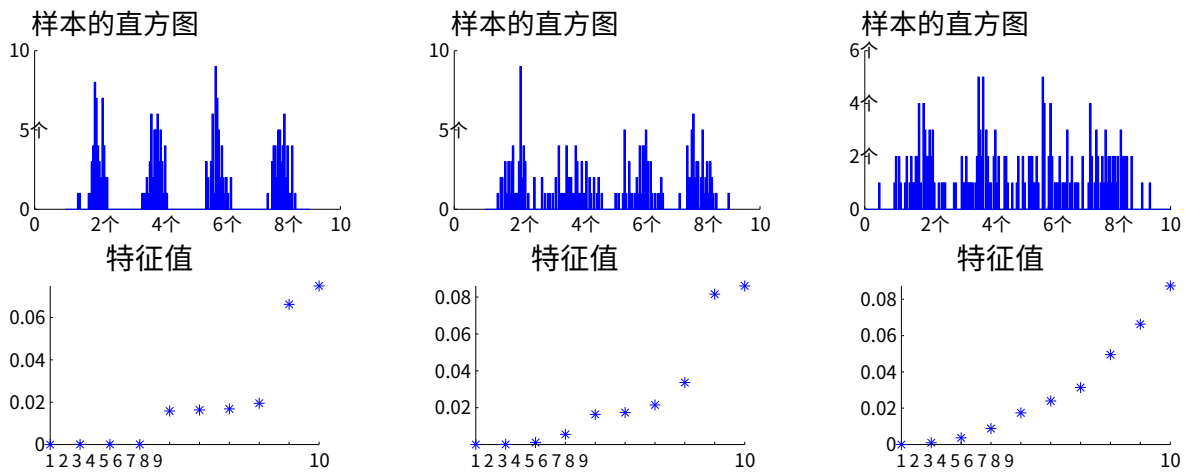


图 4：三个数据集，以及最小的 10 个特征值大小写字. 有关详细信息，请参阅文本。

布劳恩和布曼，2004 年；本大卫、冯卢克斯堡和帕尔，2006 年）。当然所有这些方法也可以用于谱聚类。此外，一种专为谱聚类设计的工具是特征间隙启发式算法，它可用于所有三种图拉普拉斯算子。这里的目标是选择号码  $k$  这样所有的特征值  $\lambda_1, \dots, \lambda_k$  很小，但是  $\lambda_{k+1}$  比较大。这个程序有几个理由。第一个基于微扰理论，我们观察到在理想情况下  $k$  完全断开的簇，特征值 0 具有多重性  $k$ ，然后  $(k+1)$  第特征值  $\lambda_{k+1} > 0$ 。谱图论可以给出其他解释。在这里，图的许多几何不变量可以在图拉普拉斯算子的第一特征值的帮助下表达或界定。特别是，切割的大小与第一特征值的大小密切相关。有关此主题的更多详细信息，请参阅 Bolla (1991)、Mohar (1997) 和 Chung (1997)。

我们想在第 4 节中介绍的玩具示例中说明本征间隙启发式算法。为此，我们考虑与第 4 节中类似的数据集，但为了改变聚类的难度，我们考虑方差递增的高斯分布。图 4 的第一行显示了三个样本的直方图。我们构建了第 4 节中描述的 10 最近邻图，并绘制了归一化拉普拉斯算子的特征值大小写字在不同的样本上（非标准化拉普拉斯算子的结果相似）。第一个数据集由四个分离良好的簇组成，我们可以看到前 4 个特征值近似为 0。然后第 4 个和第 5 个特征值之间存在间隙，即  $\lambda_5 - \lambda_4$  比较大。根据 eigengap 启发式，这个差距表明数据集包含 4 个簇。完全连接图的结果也可以观察到相同的行为（已绘制在图 1 中）。所以我们可以看到，如果数据中的聚类非常明显，则启发式方法效果很好。然而，簇的噪声越大或重叠越多，这种启发式方法的效果就越差。我们可以看到，对于聚类更“模糊”的第二个数据集，第 4 和第 5 个特征值之间仍然存在差距，但检测起来不像之前那样清晰。最后，在最后一个数据集中，没有明确定义的差距，所有特征值之间的差异大致相同。但另一方面，该数据集中的聚类重叠太多，以至于许多非参数算法将难以检测聚类，除非它们对基础模型做出强有力的假设。在此特定示例中，即使对于查看直方图的人来说，也不清楚正确的聚类数量应该是多少。这说明，作为选择聚类数量的大多数方法，如果数据包含非常明显的聚类，则特征间隙启发式算法通常效果很好，但在模棱两可的情况下，它也会返回模棱两可的结果。

最后，注意簇数的选择和邻域图的连通性参数的选择是相互影响的。例如，如果邻域图的连通性参数太小以至于该图分成，比方说， $k_0$  连接组件，然后选择  $k_0$

因为簇的数量是一个有效的选择。然而，一旦邻域图连接起来，就不清楚聚类的数量和邻域图的连通性参数是如何相互作用的。聚类数量的选择和图的连通性参数的选择本身都是难题，据我们所知，关于它们的相互作用没有什么不平凡的。

## 8.4 的 $k$ 意味着一步

我们在第 4 节中介绍的三种谱聚类算法使用  $k$  意味着作为最后一步从特征向量的实值矩阵中提取最终分区。首先，请注意，使用  $k$  表示这一步的算法。事实上，正如我们从谱聚类的各种解释中看到的那样，如果数据包含表达良好的聚类，这一步应该非常简单。例如，在理想情况下，如果簇完全分离，我们知道特征向量  $\mathbf{v}_i$  和  $\mathbf{v}_j$  是分段常数。在这种情况下，所有点  $x_i$  属于同一个集群  $C$

被精确映射到样本点是  $\mathbf{v}_i$ ，即对单位向量  $\mathbf{v}_i \in \mathbb{R}^k$ 。在这种微不足道的情况下，任何应用于点的聚类算法是  $\mathbf{v}_i \in \mathbb{R}^k$  将能够提取正确的簇。

虽然在谱聚类的最后一步中选择什么聚类算法有点随意，但可以争辩说至少点之间的欧几里得距离是  $\|\mathbf{v}_i - \mathbf{v}_j\|$  是一个有意义的数量。我们已经看到点之间的欧几里得距离是  $\|\mathbf{v}_i - \mathbf{v}_j\|$  与图表上的“通勤距离”相关，在 Nadler、Lafon、Coifman 和 Kevrekidis (2006) 中，作者表明是  $\|\mathbf{v}_i - \mathbf{v}_j\|$  也与更一般的“扩散距离”有关。此外，谱嵌入的其他用途（例如，Bolla (1991) 或 Belkin 和 Niyogi (2003)）表明，欧氏距离在  $\mathbb{R}^d$  是有意义的。

代替  $k$  意味着，人们还使用其他技术从实值表示构造他的最终解决方案。例如，在 Lang (2006) 中，作者为此目的使用了超平面。Bach 和 Jordan (2004) 提出了更高级的特征向量后处理。在这里，作者研究了第一个跨越的子空间  $k$  特征向量，并尝试使用分段常数向量尽可能好地近似这个子空间。这也导致最小化空间中的某些欧几里得距离  $R_k$ ，这可以通过一些加权来完成  $k$  意味着算法。

## 8.5 应使用哪个图拉普拉斯算子？

与谱聚类相关的一个基本问题是应该使用三个图拉普拉斯算子中的哪一个来计算特征向量。在决定这个问题之前，应该始终查看相似度图的度分布。如果图非常规则并且大多数顶点的度数大致相同，则所有拉普拉斯算子彼此非常相似，并且同样适用于聚类。但是，如果图中的度数分布非常广泛，则拉普拉斯算子会有很大差异。在我们看来，有几个论点提倡使用归一化而不是非归一化的谱聚类，并且在归一化的情况下使用特征向量  $\mathbf{v}_i$  而不是那些  $\mathbf{v}_i$  符号。

不同算法满足的聚类目标

支持归一化谱聚类的第一个论据来自图划分的观点。为简单起见，让我们讨论这个案例  $k=2$ 。一般来说，聚类有两个不同的目标：



1. 我们想要找到一个分区, 使得不同集群中的点彼此不相似, 即我们想要最小化集群间的相似性。在图形设置中, 这意味着最小化  $\text{cut}(A, B)$ 。
2. 我们想要找到一个分区, 使得同一簇中的点彼此相似, 即我们想要最大化簇内相似性  $W(A, A)$  和  $W(B, B)$ 。

RatioCut 和 Ncut 都通过显式合并  $\text{cut}(A, B)$  在目标函数中。然而, 关于第二点, 两种算法的行为不同。注意

$$W(A, A) = W(A, V) - \text{cut}(A, B) = \text{vol}(A) - \text{cut}(A, B).$$

因此, 如果  $\text{cut}(A, B)$  是小和如果体积  $\text{vol}(A)$  很大。由于这正是我们通过最小化 Ncut 实现的目标, Ncut 标准实现了第二个目标。通过考虑另一个图形切割目标函数, 即由 Ding、He、Zha、Gu 和 Simon (2001) 引入的 MinMaxCut 标准, 可以更明确地看出这一点:

$$\text{最小最大切割}(A_1, \dots, A_k) := \frac{\sum_{i=1}^k \text{cut}(A_i, V - A_i)}{\sum_{i=1}^k W(A_i, A_i)}.$$

与 Ncut 相比, 它具有  $\text{vol}(A) = \text{cut}(A, B) + W(A, A)$  在分母中, MinMaxCut 准则只有  $W(A, A)$  在分母中。在实践中, Ncut 和 MinMaxCut 通常通过类似的切割来最小化, 因为好的 Ncut 解决方案将具有较小的切割值  $\text{cut}(A, B)$  无论如何, 因此分母毕竟没有太大不同。此外, 放宽 MinMaxCut 会导致与放宽 Ncut 完全相同的优化问题, 即使用特征向量的归一化谱聚类 (Lloyd, 2006)。因此, 可以通过多种方式看出, 归一化谱聚类结合了上述两个聚类目标。

现在考虑 RatioCut 的情况。这里的目标是最大化  $\text{vol}(A)/\text{vol}(B)$  而不是音量  $\text{vol}(A)$  和体积  $\text{vol}(B)$ 。但  $\text{vol}(A)/\text{vol}(B)$  不一定与簇内相似性相关, 因为簇内相似性取决于边而不是顶点数  $n$ 。例如, 想想一个集合  $A$  它有很多顶点, 所有顶点彼此之间只有非常低的权重边。最小化 RatioCut 并不试图最大化簇内相似性, 并且对于通过非归一化光谱聚类的松弛也是如此。

因此, 这是我们要牢记的第一个重点: 归一化谱聚类实现了上述两个聚类目标, 而非归一化谱聚类仅实现了第一个目标。

## 一致性问题

对归一化谱聚类的优越性的一个完全不同的论据来自对两种算法的统计分析。在统计设置中, 假设数据点  $X_1, \dots, X_n$

已经根据某种概率分布进行了 iid 采样  $P$  在一些底层数据空间  $X$ 。那么最根本的问题就是一致性问题: 如果我们绘制越来越多的数据点, 谱聚类的聚类结果是否收敛到底层空间的一个有用分区  $X$ ?

对于两种归一化谱聚类算法, 都可以证明确实如此 (von Luxburg, Bousquet, and Belkin, 2004, 2005; von Luxburg, Bousquet, and Bousquet, to appear)。在数学上, 可以证明当我们取极限  $n \rightarrow \infty$ , 矩阵  $L$  强烈收敛

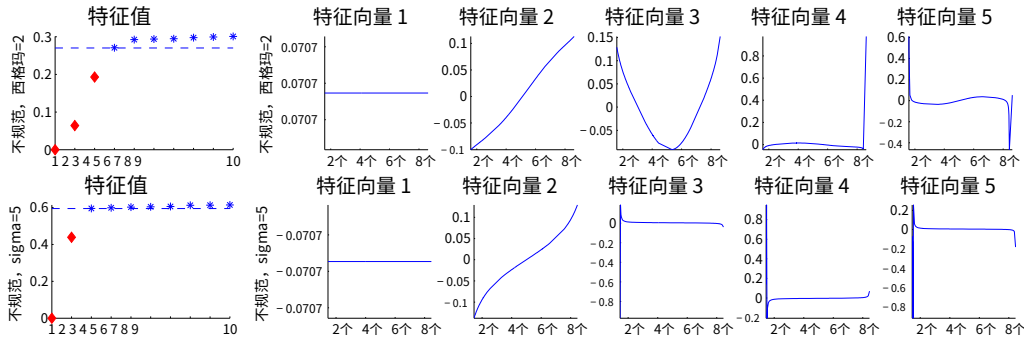


图 5: 非标准化光谱聚类的一致性。绘制的是特征值和特征向量大号, 对于参数  $\sigma=2$  (第一行) 和  $\sigma=5$  (第二行)。虚线表示最小  $d_i$ , 特征值低于  $\min d_i$  被绘制为红色菱形, 特征值高于最小值  $d_i$  被绘制为蓝色星星。有关详细信息, 请参阅文本。

给操作员  $i$  在空间上  $C(X)$  上的连续函数  $X$ . 这种收敛意味着特征值和特征向量大号符号收敛于那些  $\tilde{u}$ , 这又可以转化为关于归一化谱聚类收敛的陈述。可以证明在  $X$  通过的特征向量  $i$  可以解释为类似于谱聚类的随机游走解释。也就是说, 如果我们考虑数据空间上的扩散过程  $X$ , 然后由特征向量引起的划分  $i$  是这样的, 扩散不会经常在不同的集群之间转换 (von Luxburg et al., 2004)。关于归一化谱聚类的所有一致性声明都适用, 对于两者大号符号和大号写字, 在非常温和的条件下, 通常在现实世界的应用中得到满足。不幸的是, 解释有关这些结果的更多细节超出了本教程的范围, 因此我们建议感兴趣的读者阅读 von Luxburg 等人的文章。(出现)。

与归一化谱聚类的明确收敛性陈述相比, 非归一化谱聚类的情况要令人不快得多。可以证明, 非归一化谱聚类可能无法收敛, 或者它可以收敛到构建由数据空间的一个单点组成的聚类的平凡解 (von Luxburg et al., 2005, to appear)。在数学上, 即使可以证明矩阵  $(1/n)$  大号本身收敛于某个极限算子  $\mathcal{L}$  上  $C(X)$  作为  $n \rightarrow \infty$ , 这个极限算子的光谱性质  $\mathcal{L}$  可能非常讨厌, 以至于它们阻止了谱聚类的收敛。可以构造示例来表明这不仅是非常大的样本量的问题, 而且即使对于小样本量也可能导致完全不可靠的结果。至少可以描述那些问题不发生时的条件: 我们必须确保特征值大号对应于非归一化谱聚类中使用的特征向量的度数明显小于图中的最小度数。这意味着如果我们使用第一个  $k$  用于聚类的特征向量, 然后  $\lambda_{k+1}$

分钟  $j=1, \dots, n d_j$  应该适用于所有人  $i=1, \dots, k$ .

这种情况的数学原因是对应于大于  $\min$  的特征值的特征向量  $d_j$  近似 Dirac 函数, 即它们在除一个坐标外的所有坐标中都近似为 0。如果这些特征向量用于聚类, 那么它们将特征向量不为零的一个顶点与所有其他顶点分开, 我们显然不想构造这样的分区。我们再次参考文献以获得精确的陈述和证明。

为了说明这种现象, 再次考虑第 4 节中的玩具数据集。我们考虑基于全连接图的非归一化图拉普拉斯算子的第一个特征值和特征向量, 用于不同的参数选择  $\sigma$  高斯相似函数 (参见图 1 的最后一行和图 5 的所有行)。特征值高于最小值  $d_i$  被绘制为蓝色星星, 特征值

低于最低限度  $d_i$  是  $p$  被绘制为红色钻石。虚线表示最小  $d_i$ . 一般来说, 我们可以看到

与远低于虚线的特征值对应的特征向量是“有用的”特征向量。如果  $\sigma = 1$ （已经绘制在图 1 的最后一行），特征值 2、3 和 4 明显低于最小值  $d_j$ ，并且相应的特征向量 2、3 和 4 是有意义的（如第 4 节所述）。如果我们增加参数  $\sigma$ ，我们可以观察到特征值趋向于向最小值移动  $d_j$ 。如果  $\sigma = 2$ ，只有前三个特征值低于最小值  $d_j$ （见图 5 中的第一行），以防万一  $\sigma = 5$  只有前两个特征值低于最小值  $d_j$ （图 5 中的第二行）。我们可以看到，一旦特征值接近或高于最小值  $d_j$ ，其对应的特征向量近似狄拉克函数。当然，那些特征向量不适合构建聚类。在极限为  $n \rightarrow \infty$ ，这些特征向量会收敛到完美的狄拉克函数。我们对有限样本案例的说明表明，这种行为不仅发生在大样本量的情况下，而且甚至可以在我们的玩具数据集中的小样本上产生。

非常重要是要强调这些问题只涉及矩阵的特征向量，它们不会出现  $\frac{1}{\sqrt{n}}$  或者  $\frac{1}{\sqrt{m}}$ 。因此，从统计的角度来看，最好避免非归一化的谱聚类，而改用归一化的算法。

哪个归一化拉普拉斯算子？

查看两种归一化谱聚类算法之间的差异  $\frac{1}{\sqrt{n}}$  和  $\frac{1}{\sqrt{m}}$ ，谱聚类的所有三种解释都支持  $\frac{1}{\sqrt{n}}$ 。原因是特征向量  $\frac{1}{\sqrt{n}}$  是聚类指标向量  $\frac{1}{\sqrt{n}}$ ，而特征向量  $\frac{1}{\sqrt{m}}$  另外乘以  $\sqrt{\frac{n}{m}}$ ，这可能会导致不需要的伪像。作为使用  $\frac{1}{\sqrt{m}}$  也没有任何计算优势，因此我们提倡使用  $\frac{1}{\sqrt{n}}$ 。

## 9 展望与延伸阅读

谱聚类可以追溯到 Donath 和 Hoffman (1973)，他们首先建议基于邻接矩阵的特征向量构建图分区。同年，Fiedler (1973) 发现图的二分与拉普拉斯图的第二个特征向量密切相关，他建议用这个特征向量对图进行分割。从那时起，谱聚类已在不同社区中多次被发现、重新发现和扩展，例如参见 Pothén、Simon 和 Liou (1990)、Simon (1991)、Bolla (1991)、Hagen 和 Kahng (1992)、Hendrickson 和 Leland (1995)、Van Driessche 和 Roose (1995)、Barnard、Pothén 和 Simon (1995)、Spielman 和 Teng (1996)、Guattery 和 Miller (1998)。在 Spielman 和 Teng (1996) 中可以找到对谱聚类历史的一个很好的概述。

在机器学习社区中，谱聚类已因 Shi 和 Malik (2000)、Ng 等人的作品而流行起来。(2002)，Meila 和 Shi (2001)，以及 Ding (2004)。随后，谱聚类已扩展到许多非标准设置，例如应用于共聚类问题的谱聚类 (Dhillon, 2001)，具有附加边信息的谱聚类 (Joachims, 2003) 谱聚类与加权内核之间的连接- $k$ -均值算法 (Dhillon、Guan 和 Kulis, 2005 年)、基于谱聚类的学习相似函数 (Bach 和 Jordan, 2004 年) 或分布式环境中的谱聚类 (Kempe 和 McSherry, 2004 年)。此外，还发现了关于谱聚类与其他算法之间关系的新理论见解。谱聚类和加权核之间的联系  $k$ -均值算法在 Dhillon 等人中有描述。(2005)。谱聚类和 (核) 主成分分析之间的关系依赖于这样一个事实，即图拉普拉斯算子的最小特征向量也可以解释为核矩阵 (Gram 矩阵) 的最大特征向量。这种解释有两种不同的风格：而 Bengio 等人。(2004) 解释矩阵  $\frac{1}{\sqrt{n}}$  作为内核矩阵，其他作者 (Saerens,

Fouss、Yen 和 Dupont, 2004 年) 解释了 Moore-Penrose 逆函数  $L^+$  或者  $L^\dagger$  符号作为核矩阵。这两种解释都可用于构造 (不同的) 样本外扩展以进行谱聚类。关于谱聚类的应用案例, 近几年在各个科学领域发表的论文数量庞大, 无法一一列举。我们鼓励读者使用 “谱聚类” 一词查询他最喜欢的文献数据库, 以了解应用程序的多样性。

谱聚类的成功主要在于它没有对聚类的形式做出强假设。相对于  $k$ -意味着, 在生成的聚类形成凸集 (或者更准确地说, 位于底层空间的不相交凸集) 的情况下, 谱聚类可以解决非常普遍的问题, 例如螺旋交织。此外, 只要我们确保相似图是稀疏的, 即使对于大型数据集也可以有效地实现谱聚类。一旦选择了相似图, 我们只需要解决一个线性问题, 不存在陷入局部极小或多次重新启动算法并使用不同初始化的问题。然而, 我们已经提到选择一个好的相似度图并非易事, 并且在邻域图参数选择不同的情况下, 谱聚类可能会非常不稳定。因此, 谱聚类不能充当 “黑盒算法”, 自动检测任何给定数据集中的正确聚类。但它可以被认为是一个强大的工具, 如果小心使用, 可以产生良好的效果。

在机器学习领域, 图拉普拉斯算子不仅用于聚类, 而且还用于许多其他任务, 例如半监督学习 (例如, Chapelle、Schölkopf 和 Zien, 2006 年的概述) 或流形重建 (例如, 贝尔金和 Niyogi, 2003 年)。在大多数应用中, 图拉普拉斯算子用于对 “接近” 的数据点 (即,  $w_{ij}$  大) 应该有一个 “相似” 的标签 (即  $F_i \approx F_j$ )。一个功能  $F$  满足这个假设, 如果  $w_{ij}(F_i - F_j)^2 \uparrow$

对所有人来说都很小我,  $j$ , 那是  $F$  如果是小。凭借这种直觉, 我们可以使用二次形式  $F$  如果作为转导分类问题中的正则化项。另一种解释图拉普拉斯算子使用的方法是通过它们编码的平滑假设。一个功能  $F$  它的价值很低  $F$  如果具有这样的特性, 即它在数据点密集的区域 (即图形紧密连接) 仅 “一点点” 变化, 而在低数据区域允许变化更大 (例如, 改变符号) 密度。从这个意义上说, 一个小值  $F$  如果在半监督学习中编码所谓的 “集群假设”, 它要求分类器的决策边界应位于低密度区域。

一个经常使用的直觉是, 图拉普拉斯算子形式上看起来像一个连续的拉普拉斯算子 (这也是 “图拉普拉斯算子” 这个名字的来源)。要看到这一点, 请转换局部相似性  $w_{ij}$  到远处  $d_{ij}$  通过关系  $w_{ij} = 1/d_{ij}^2$  并观察到

$$w_{ij}(F_i - F_j)^2 \approx \frac{(F_i - F_j)^2}{d_{ij}^2}$$

看起来像一个差商。因此, 等式  $F$  如果=  $\sum_{ij} w_{ij}(F_i - F_j)^2$  从命题 1 看起来像是与标准拉普拉斯算子相关联的二次型的离散版本大号上  $R_n$ , 满足

$$\int \langle G, \Delta u \rangle = \int |\nabla u|^2 dx_0$$

这种直觉在 Belkin (2003)、Lafon (2004)、Hein、Audibert 和 von Luxburg (2005) 的作品中得到了精确体现; M.、Audibert 和 von Luxburg (2007)、Belkin 和 Niyogi (2005)、Hein (2006)、Giné 和 Koltchinskii (2005)。一般来说, 证明图拉普拉斯算子是某些连续拉普拉斯算子的离散版本, 并且如果拉普拉斯图是在随机采样数据点的相似图上构建的, 则它收敛到某个连续拉普拉斯算子 (或

Laplace-Beltrami 算子) 在底层空间上。Belkin (2003) 研究了收敛证明的第一个重要步骤, 它处理与离散图 Laplacians 相关的连续算子收敛到 Laplace-Beltrami 算子。Lafon (2004) 将他的结果从均匀分布推广到一般分布。然后在 Belkin 和 Niyogi (2005) 中, 作者在具有均匀分布的流形上使用高斯相似函数证明了非归一化拉普拉斯算子的逐点收敛结果。与此同时, 海因等人。(2005) 证明更一般的结果, 考虑到所有不同的拉普拉斯图大小, 大小写字, 和大小符号, 更一般的相似函数, 以及具有任意分布的流形。在 Giné 和 Koltchinskii (2005) 中, 分布和均匀收敛结果在具有均匀分布的流形上得到证明。Hein (2006) 研究了由图 Laplacians 引起的平滑泛函的收敛, 并显示了一致的收敛结果。

除了将图拉普拉斯算子应用于最广泛意义上的划分问题外, 图拉普拉斯算子还可以用于完全不同的目的, 例如图形绘制 (Koren, 2005)。事实上, 图的拓扑和属性与图拉普拉斯矩阵之间的联系比我们在本教程中提到的要紧密得多。现在对最基本的属性有了了解, 有兴趣的读者可以自己进一步探索和欣赏该领域的大量文献。

## 参考

- Aldous, D. 和 Fill, J. (准备中)。可逆马尔可夫链和图上的随机游走。  
在线版本可在 <http://www.stat.berkeley.edu/users/aldous/RWG/book.html> 获得。Bach, F. 和 Jordan, M. (2004)。学习谱聚类。在 S. Thrun、L. Saul 和 B. Schölkopf (编辑), *神经信息处理系统进展 16 (NIPS)*(第 305 – 312 页)。马萨诸塞州剑桥市: 麻省理工学院出版社。
- Bapat, R.、Gutman, I. 和 Xiao, W. (2003)。一种计算阻力距离的简单方法。Z. *自然资源*, 58, 494 – 498。
- Barnard, S.、Pothen, A. 和 Simon, H. (1995)。一种稀疏包络缩减的谱算法  
矩阵。数值线性代数及其应用, 2(4), 317 – 334。
- 贝尔金, M. (2003)。学习流形的问题。芝加哥大学博士论文。Belkin, M. 和 Niyogi, P. (2003)。用于降维和数据表示的拉普拉斯特征图  
发送。神经计算, 15(6), 1373 – 1396。
- Belkin, M. 和 Niyogi, P. (2005)。迈向基于拉普拉斯流形的理论基础  
方法。在 P. Auer 和 R. Meir (编辑) 中, 第 18 届学习理论年会 (COLT) 论文集(第 486 – 500 页)。斯普林格, 纽约。
- Ben-David, S.、von Luxburg, U. 和 Pál, D. (2006)。对聚类稳定性的清醒认识。在 G. Lugosi 和 H. Simon (编辑), 第 19 届学习理论年会 (COLT) 论文集 (第 5-19 页)。施普林格, 柏林。
- Bengio, Y.、Delalleau, O.、Roux, N.、Paiement, J.、Vincent, P. 和 Ouimet, M. (2004 年)。学习  
特征函数链接谱嵌入和内核 PCA。神经计算, 16, 2197 – 2219。Ben-Hur, A.、Elisseeff, A. 和 Guyon, I. (2002)。一种基于稳定性的结构发现方法  
在聚类数据中。在太平洋生物计算研讨会(第 6 – 17 页)。Bhatia, R. (1997)。矩阵分析。斯普林格, 纽约。
- Bie, TD 和 Cristianini, N. (2006)。图形切割聚类、转换和转换的快速 SDP 松弛  
其他组合问题。JMLR, 7, 1409 – 1436。
- Bolla, M. (1991)。多图谱和分类属性之间的关系(技术再  
端口号 DIMACS-91-27)。离散数学和理论计算机科学中心。Brémaud, P. (1999)。马尔可夫链: 吉布斯域、蒙特卡罗模拟和队列。纽约: 施普林格出版社。

- Brito, M.、Chavez, E.、Quiroz, A. 和 Yukich, J. (1997)。相互  $k$  最近邻的连通性  
聚类和异常值检测中的图形。《统计和概率快报》, 35, 33 – 42. Bui, TN 和 Jones, C. (1992)。找到好的近似顶点和边缘分区是 NP-hard。  
信息。过程。莱特。 , 42(3), 153 – 159。
- Chapelle, O.、Schölkopf, B. 和 Zien, A. (编辑)。(2006)。《半监督学习》。麻省理工学院出版社, 剑桥。
- 钟 F. (1997)。《谱图论卷》。数学方面的 CBMS 区域会议系列的 92-  
数学)。华盛顿数学科学会议委员会。
- Dhillon, I. (2001)。使用二分谱图分区对文档和单词进行联合聚类。在  
《第七届 ACM SIGKDD 知识发现与数据挖掘 (KDD) 国际会议论文集》(第 269 – 274 页)。纽约: ACM 出版社。
- Dhillon, I.、Guan, Y. 和 Kulis, B. (2005)。《统一的内核视图  $k$ -意味着, 谱聚类, 和  
图划分》(技术报告编号 UTCS TR-04-25)。德克萨斯大学奥斯汀分校。丁, C. (2004 年)。《谱聚  
类教程》。在 ICML 上发表的演讲。(幻灯片可在  
<http://crd.lbl.gov/-光盘/光谱/>) Ding, C.、He,  
X.、Zha, H.、Gu, M. 和 Simon, H. (2001)。图的最小-最大切割算法  
分区和数据聚类。在《第一届 IEEE 数据挖掘国际会议 (ICDM) 论文集》(第 107 – 114 页)。美国华盛顿  
特区: IEEE 计算机协会。Donath, WE 和 Hoffman, AJ (1973)。图划分的下界。《IBM J. Res.  
开发》, 17, 420 – 425。
- M. 菲德勒 (1973)。图的代数连通性。《捷克斯洛伐克数学》, 23, 298 – 305. Fouss, F.、Pirrotte, A.、  
Renders, J.-M. 和 Saerens, M. (2007)。相似的随机游走计算  
图的节点之间的关系在协作推荐中的应用。《IEEE 跨。知识 数据工程》, 19(3), 355–369。
- Fraley, C. 和 Raftery, AE (2002)。基于模型的聚类、判别分析和密度  
估计。《日本航空航天局》, 97, 611 – 631。
- Giné, E. 和 Koltchinskii, V. (2005)。Laplace-Beltrami 的经验图拉普拉斯近似  
运营商: 大样本结果。在《第四届高维概率国际会议论文集》(第 238-259 页)。
- Golub, G. 和 Van Loan, C. (1996)。《矩阵计算》。巴尔的摩: 约翰霍普金斯大学  
按。
- Guattery, S. 和 Miller, G. (1998)。关于光谱分离器的质量。《SIAM 矩阵杂志  
肛门。申请》, 19(3), 701 – 719。
- Gutman, I. 和 Xiao, W. (2004)。拉普拉斯矩阵的广义逆和一些应用。  
《Bulletin de l'Academie Serbe des Sciences at des Arts (Cl. Math. Natur.)》, 129, 15 – 23. Hagen,  
L. 和 Kahng, A. (1992)。用于比率切割分区和聚类的新谱方法。  
《IEEE 跨。计算机辅助设计》, 11(9), 1074 – 1085。
- Hastie, T.、Tibshirani, R. 和 Friedman, J. (2001)。《统计学习的要素》。纽约:  
施普林格。
- 海因, M. (2006 年)。基于自适应图的正则化的均匀收敛。在《程序的  
第 19 届学习理论年会 (COLT)》(第 50 – 64 页)。斯普林格, 纽约。Hein, M.、Audibert, J.-Y. 和 von  
Luxburg, U. (2005)。从图到流形——弱与强  
图拉普拉斯算子的逐点一致性。在 P. Auer 和 R. Meir (编辑) 中, 《第 18 届学习理论年会 (COLT) 论  
文集》(第 470 – 485 页)。斯普林格, 纽约。Hendrickson, B. 和 Leland, R. (1995)。一种改进的映射  
谱图划分算法  
并行计算。《SIAM J. 论科学计算》, 16, 452 – 469。
- 约阿希姆斯 (2003)。通过谱图划分的转换学习。在 T. Fawcett 和  
N. Mishra (编辑), 《第 20 届机器学习国际会议 (ICML) 论文集》(第 290 – 297 页)。AAAI 出版社。
- Kannan, R.、Vempala, S. 和 Vetta, A. (2004)。关于聚类: 好、坏和光谱。《杂志  
美国计算机协会》, 51(3), 497–515。
- Kempe, D. 和 McSherry, F. (2004)。一种用于光谱分析的分散算法。在《诉讼程序

第 36 届 ACM 计算理论研讨会 (STOC)(第 561 – 568 页)。美国纽约州纽约市: ACM 出版社。

- Klein, D. 和 Randic, M. (1993)。阻力距离。数理化学学报,12, 81 – 95. Koren, Y. (2005)。通过特征向量绘制图形: 理论与实践。计算机和数学与应用程序,49, 1867 – 1888.
- Lafon, S. (2004)。扩散图和几何谐波。耶鲁大学博士论文。
- Lang, K. (2006)。修复光谱方法的两个弱点。在 Y. Weiss、B. Schölkopf 和 J. Platt 中 (编辑), 神经信息处理系统的进展 18(第 715 – 722 页)。马萨诸塞州剑桥市: 麻省理工学院出版社。
- Lange, T.、Roth, V.、Braun, M. 和 Buhmann, J. (2004 年)。基于稳定性的聚类验证解决方案。神经计算,16(6), 1299 – 1323。
- Lovász, L. (1993)。图上的随机游走: 一项调查。在组合数学, Paul Erdős 八十岁(页数 353 – 397)。布达佩斯: János Bolyai Math. 社会。
- Lütkepohl, H. (1997)。矩阵手册。奇切斯特: 威利。
- M.、Audibert, J.-Y. 和 von Luxburg, U. (2007 年)。图拉普拉斯算子及其在随机邻域图上的收敛。JMLR, 8个, 1325 – 1370。
- Meila, M. 和 Shi, J. (2001)。光谱分割的随机游走视图。在第八届国际人工智能与统计研讨会 (AISTATS)。
- Mohar, B. (1991)。图的拉普拉斯谱。在图论、组合学和应用。卷。2 (卡拉马祖, 密歇根州, 1988 年) (第 871 – 898 页)。纽约: 威利。
- Mohar, B. (1997)。图的拉普拉斯特征值的一些应用。在 G. Hahn 和 G. Sabidussi (编辑), 图对称性: 代数方法与应用(卷。北约 ASI 系列 C 497, 第 225-275 页)。克鲁沃。
- Nadler, B.、Lafon, S.、Coifman, R. 和 Kevrekidis, I. (2006)。扩散图、光谱聚类 和 Fokker-Planck 算子的特征函数。在 Y. Weiss、B. Schölkopf 和 J. Platt (编辑) 中, 神经信息处理系统的进展 18(第 955 – 962 页)。马萨诸塞州剑桥市: 麻省理工学院出版社。
- Ng, A.、Jordan, M. 和 Weiss, Y. (2002)。关于谱聚类: 分析和算法。在 T. Dietterich、S. Becker 和 Z. Ghahramani (编辑), 神经信息处理系统的进展 14(第 849 – 856 页)。麻省理工学院出版社。
- J. 诺里斯 (1997)。马尔可夫链。剑桥: 剑桥大学出版社。
- 彭罗斯, M. (1999)。最小生成树最长边的强定律。安。的概率。 , 27(1), 246 – 260。
- Pothen, A.、Simon, HD 和 Liou, KP (1990)。用特征向量分割稀疏矩阵图。SIAM 矩阵肛门杂志。申请,11, 430 – 452。
- Saerens, M.、Fouss, F.、Yen, L. 和 Dupont, P. (2004)。图的主成分分析, 及其与谱聚类的关系。在第 15 届欧洲机器学习会议 (ECML) 论文集(第 371 – 383 页)。施普林格, 柏林。
- Shi, J. 和 Malik, J. (2000)。归一化切割和图像分割。IEEE 模式汇刊分析和机器学习,22(8), 888 – 905。
- 西蒙, H. (1991)。为并行处理划分非结构化问题。计算系统工程,2个, 135 – 148。
- Spielman, D. 和 Teng, S. (1996)。谱划分工作: 平面图和有限元网格。在第 37 届计算机科学基础年会 (佛蒙特州伯灵顿, 1996 年) (第 96 – 105 页)。加利福尼亚州洛斯阿拉米托斯: IEEE 计算。社会。按。(另请参阅扩展技术报告。)
- Stewart, G. 和 Sun, J. (1990)。矩阵微扰理论。纽约: 学术出版社。
- Still, S. 和 Bialek, W. (2004)。有多少簇? 信息论的观点。神经的电脑。 ,16(12), 2483 – 2506。
- Stoer, M. 和 Wagner, F. (1997)。一个简单的最小割算法。J.ACM,44(4), 585 – 591. Tibshirani, R.、Walther, G. 和 Hastie, T. (2001)。通过以下方式估计数据集中的簇数差距统计。J.皇家。国家主义者。社会。乙,63(2), 411 – 423。

- Van Driessche, R. 和 Roose, D. (1995)。一种改进的谱二分法及其应用  
到动态负载均衡。 *并行计算*, 21(1), 29 – 48。
- von Luxburg, U.、Belkin, M. 和 Bousquet, O. (出场)。谱聚类的一致性。 *志  
统计学*。 (另见技术报告 134, 马克斯普朗克生物控制论研究所, 2004 年)
- von Luxburg, U.、Bousquet, O. 和 Belkin, M. (2004)。关于谱聚类的收敛性  
随机样本：归一化的情况。在 J. Shawe-Taylor 和 Y. Singer (编辑) 中, *第 17 届学习理论年会  
(COLT) 论文集*(第 457 – 471 页)。斯普林格, 纽约。 von Luxburg, U.、Bousquet, O. 和 Belkin, M.  
(2005)。谱聚类的局限性。在 L. Saul 中,  
Y. Weiss 和 L. Bottou (编辑), *神经信息处理系统 (NIPS) 的进展 17* (第 857 – 864 页)。马萨诸塞  
州剑桥市：麻省理工学院出版社。
- Wagner, D. 和 Wagner, F. (1993)。在最小割和图二分之间。在 *程序的  
第 18 届计算机科学数学基础国际研讨会 (MFCS)*(第 744 – 750 页)。伦敦：施普林格。