

文章编号: 1001-0920(2011)01-0080-05

## 单位超球面上的二元聚类算法

程 成, 萧蕴诗, 岳继光

(同济大学 电子与信息工程学院, 上海 201804)

**摘 要:** 鉴于在数据学习理论中, 基于各种几何正则化的处理方法引起人们的广泛关注, 以基于独立成分分析的地震数据处理为背景, 针对数据向量的单位模长约束, 研究了单位超球面上的二元聚类问题. 通过欧氏诱导度量, 推导了单位超球面上的黎曼梯度公式, 并据此构造了求取其上二元聚类和数据平均的不动点迭代算法. 实验结果表明了其有效性和优越性.

**关键词:** 单位超球面; 聚类; 黎曼梯度; 测地距离

中图分类号: TN911; O186

文献标识码: A

## Clustering algorithm for two classes on unit Hypersphere

CHENG Cheng, XIAO Yun-shi, YUE Ji-guang

(School of Electronics and Information Engineering, Tongji University, Shanghai 201804, China. Correspondent: CHENG Cheng, E-mail: chengcheng\_lcc@163.com)

**Abstract:** The methods of geometrical regularization in data learning theory have caused wide public concern. In the background of seismic data processing based on independent component analysis, a problem of two-class clustering on unit hypersphere is studied for the unit norm constraint. The Riemannian gradient is formulated based on the induced metric from Euclidean space, which realizes the construction of a fixed-point algorithm for two-class clustering and data averaging on unit hypersphere. Finally, the simulation result shows the effectiveness and superiority of this algorithm.

**Key words:** unit hypersphere; clustering; riemannian gradient; geodesic distance

### 1 引 言

自经典的流形学习算法局部线性嵌入 (LLE)<sup>[1]</sup>和等距映射 (ISOMAP)<sup>[2]</sup>提出以来, 基于几何正则化的数据处理方法开始引起人们的广泛关注<sup>[3]</sup>. 流形学习算法以局部度量代替数据空间的整体欧氏度量, 适应了数据空间弯曲的内蕴特征. 在流形上, 连接任意两点间的最短路径称为测地线, 而其长度称为测地距离, 基于测地距离的数据聚类算法是通常  $K$ -means 算法在弯曲流形上的自然推广. Goh 等人<sup>[4]</sup>通过二次均方重整化对参数空间赋予了特殊的度量结构, 并据此研究了概率密度函数簇在黎曼流形上的聚类问题, 这为数据空间的概率密度估计和统计回归分析提供了一种新的方法, 并将此方法应用于图像分割, 取得了满意的结果<sup>[5]</sup>. 文献<sup>[6]</sup>提出了一种称为软测地核  $K$ -means 聚类的算法, 该方法利用测地距离构造核函数, 保证了数据集在特征空间的线性可分性, 也克服了通常方法中对核函数选取的随意性.

在数据挖掘中, 相对于主成分分析 (PCA), 基于独立成分分析 (ICA) 的特征提取能够获得更有意义的结果<sup>[7-8]</sup>. 本文基于 ICA 研究了地震数据处理中的子波提取问题. 在提取过程中, 为了保证算法的稳定性, 对特征向量, 即地震子波施加的单位范数约束条件赋予了特征空间一个自然的单位超球面几何结构. 又由于 ICA 中符号的不确定性, 使得提取出的特征向量呈现出近似对称的二元分布特点. 对此, 本文通过构造黎曼度量, 研究了单位超球面上的二元聚类分析, 并提出了相应的不动点迭代算法. 在特征维数上, 单位超球面与其嵌入的欧氏空间相差 1, 因此这种度量适应了空间弯曲的结构特征. 鉴于在子波提取中数据向量近似对称的分布特点, 将聚类结果中的一类数据反号, 就可以自然地得到求取这些数据向量单位超球面平均的算法. 该方法不仅能够平滑数据中的不规则点, 还可以消除测量系统中的随机影响. 最后, 将该二元聚类分析和数据平均的方法应用于多道模拟地震

收稿日期: 2009-10-29; 修回日期: 2010-02-02.

基金项目: 国家自然科学基金项目(40872090).

作者简介: 程成(1980—), 男, 博士生, 从事盲源分离、机器学习等研究; 萧蕴诗(1946—), 男, 教授, 博士生导师, 从事信号处理、智能控制理论等研究.

数据的子波提取, 经比较发现其结果明显优于基于欧氏度量的  $K$ -means 聚类所得的结果.

## 2 问题描述

假设欧氏空间  $\mathbf{R}^p$  中一数据点集记为  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ ,  $K$ -means 聚类即是将这  $N$  个数据点分别划归到  $K$  个不同的类别中, 其中类别数  $K$  是预先已知的, 而每个类别分别以聚类中心  $\mathbf{c}_k$  标记. 任意数据点  $\mathbf{x}_n$  与各聚类中心  $\mathbf{c}_k$  间距离的大小作为其划归相应类别的依据. 在  $\mathbf{R}^p$  中,  $K$ -means 聚类问题可以描述为如下优化问题:

$$\arg \min_{r_{nk}, \mathbf{c}_k} J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \mathbf{c}_k\|^2, \quad (1)$$

其中

$$r_{nk} = \begin{cases} 1, & k = \arg \min_j \|\mathbf{x}_n - \mathbf{c}_j\|^2; \\ 0, & \text{otherwise.} \end{cases}$$

$r_{nk}$  是判断数据点类别属性的二元指标, 目标函数  $J$  的建立依赖于  $\mathbf{R}^p$  中由 2-范数结构诱导的直线距离.

流形学习的结果对数据点的度量结构具有明显的依赖关系, 因此对于特殊的数据流形, 必须考虑与其相应的几何与度量结构. 假设  $d$  就是满足这种结构的二元泛函, 则相应的聚类问题可以描述为

$$\arg \min_{r_{nk}, \mathbf{c}_k} J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} d(\mathbf{x}_n, \mathbf{c}_k), \quad (2)$$

其中

$$r_{nk} = \begin{cases} 1, & k = \arg \min_j d(\mathbf{x}_n, \mathbf{c}_j); \\ 0, & \text{otherwise.} \end{cases}$$

当  $K = 2$ , 且所有参数向量具有单位模长时, 问题 (2) 就是本文所要研究的单位超球面上的二元聚类问题.

以下通过研究单位超球面上的几何结构, 建立一种自然的度量关系, 并且用于数据二元聚类和求取平均的研究.

## 3 单位超球面上的几何结构

### 3.1 数学定义

欧氏空间  $\mathbf{R}^p$  中的单位超球面可以表示为

$$\mathbf{S}^{p-1} = \{\mathbf{x} \in \mathbf{R}^p | \mathbf{x}^T \mathbf{x} = 1\}.$$

利用  $\mathbf{R}^p$  中的标准内积, 任意  $\mathbf{x} \in \mathbf{S}^{p-1}$  处的仿射空间可以正交分解为切空间和法空间的直和, 分别为

$$T_{\mathbf{x}}\mathbf{S}^{p-1} = \{\mathbf{v} \in \mathbf{R}^p | \mathbf{v}^T \mathbf{x} = 0\}, \quad N_{\mathbf{x}}\mathbf{S}^{p-1} = \{\alpha \mathbf{x} | \alpha \in \mathbf{R}\}.$$

作为  $\mathbf{R}^p$  中的一个嵌入子流形, 可以对  $\mathbf{x}$  处的切空间  $T_{\mathbf{x}}\mathbf{S}^{p-1}$  赋予一个诱导度量, 使其局部欧氏化. 该诱导度量为

$$\langle \mathbf{v}_1, \mathbf{v}_2 \rangle_{\mathbf{x}}^{\mathbf{S}^{p-1}} = \langle \mathbf{v}_1, \mathbf{v}_2 \rangle_{\mathbf{x}}^{\mathbf{R}^p} = \mathbf{v}_1^T \mathbf{v}_2. \quad (3)$$

### 3.2 黎曼梯度

设  $f$  是  $\mathbf{S}^{p-1}$  上的一个光滑函数. 作为嵌入子流

形,  $f$  也可以看作是  $\mathbf{R}^p$  中的函数. 根据微积分理论,  $f$  在  $\mathbf{R}^p$  中的梯度向量  $\partial f$  指向函数增长最快的方向. 由于  $\mathbf{S}^{p-1}$  是一个约束空间, 可以忽略  $f$  在法向上的变化而只考虑其在切空间中的梯度分量, 这就引出了黎曼梯度的概念.

根据诱导度量 (3),  $f$  所对应的黎曼梯度  $\nabla f \in T_{\mathbf{x}}\mathbf{S}^{p-1}$  定义为满足如下约束的切向量:

$$\langle \mathbf{v}, \nabla f \rangle_{\mathbf{x}}^{\mathbf{S}^{p-1}} = \langle \mathbf{v}, \partial f \rangle_{\mathbf{x}}^{\mathbf{R}^p} = \mathbf{v}^T \partial f. \quad (4)$$

其中:  $\partial f$  是  $\mathbf{R}^p$  中的常规梯度,  $\mathbf{v}$  是  $T_{\mathbf{x}}\mathbf{S}^{p-1}$  中任意一个切向量.

由于  $T_{\mathbf{x}}\mathbf{S}^{p-1}$  被赋予与  $\mathbf{R}^p$  相同的内积形式, 由式 (4) 可得  $\langle \mathbf{v}, \nabla f - \partial f \rangle = 0$ , 即存在  $\alpha \in \mathbf{R}$  使得

$$\nabla f - \partial f = \alpha \mathbf{x}. \quad (5)$$

进一步计算可知  $\alpha = -\langle \mathbf{x}, \partial f \rangle = -\mathbf{x}^T \partial f$ . 将该结果代入式 (5) 即可得到  $f$  在  $\mathbf{x}$  点处的黎曼梯度公式

$$\nabla f = (\mathbf{I}_p - \mathbf{x}\mathbf{x}^T) \partial f, \quad (6)$$

其中  $\mathbf{I}_p$  代表  $p$  阶单位矩阵.

根据诱导度量, 可以证明式 (6) 中的黎曼梯度  $\nabla f$  就是常规梯度  $\partial f$  在切空间  $T_{\mathbf{x}}\mathbf{S}^{p-1}$  中的投影向量. 以球面  $\mathbf{S}^2$  为例, 图 1 给出了各类梯度的图示说明, 其中  $\nabla^\perp f = \partial f - \nabla f$  是法向梯度.

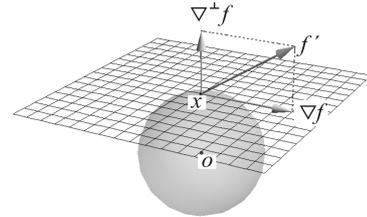


图 1 二维球面上的各种梯度向量

### 3.3 单位超球面上的测地距离

由文献 [9] 可知, 被赋予诱导度量后,  $\mathbf{S}^{p-1}$  上的测地线就是其表面的大圆弧. 通过将测地线看成质点在  $\mathbf{S}^{p-1}$  上运动的轨迹可以给出其动力学解释. 假设质点以  $\mathbf{x}$  为起点, 以  $\mathbf{v} \in T_{\mathbf{x}}\mathbf{S}^{p-1}$  为初始速度运动, 则其运动轨迹就是由这两个向量所决定的单位圆周, 测地线方程为

$$g(s) = \cos(s)\mathbf{x} + \sin(s)\mathbf{v}. \quad (7)$$

其中:  $\mathbf{x}$  和  $\mathbf{v}$  均是单位向量;  $s$  是弧长参数, 由于圆周具有单位半径,  $s$  也是弧度参数. 此处限定质点具有单位速度并非本质的约束, 因为任何曲线通过弧长参数化都具有单位速度, 并且经过参数替换, 弧长的大小是不会发生改变的 [9]. 进一步计算可知, 该测地线方程可以看成是质点在  $\mathbf{S}^{p-1}$  上做匀速圆周运动的轨迹.

以下考虑  $\mathbf{S}^{p-1}$  上任意两点  $\mathbf{x}_1$  和  $\mathbf{x}_2$  间的测地距离  $d(\mathbf{x}_1, \mathbf{x}_2)$ . 设测地线  $g(s)$  满足  $g(0) = \mathbf{x}_1, g(S) = \mathbf{x}_2$ .

由于曲线已经被弧长参数化, 连接该两点间的测地距离即为  $d(\mathbf{x}_1, \mathbf{x}_2) = S$ . 代入边界条件及单位模长约束可得  $\cos(S) = \mathbf{x}_1^T \mathbf{x}_2$ . 由于大圆的优弧与劣弧都满足测地线方程, 显然应该选取劣弧的长度作为测地距离, 而反三角余弦函数能够自动取值于劣弧范围. 经计算,  $\mathbf{S}^{p-1}$  上任意两点间的测地距离为

$$d(\mathbf{x}_1, \mathbf{x}_2) = \arccos(\mathbf{x}_1^T \mathbf{x}_2).$$

可以发现以上结果同  $\mathbf{S}^1$  的情形是一致的.

#### 4 单位超球面上的二元聚类分析

根据式 (2) 对于广义聚类问题的描述, 选取  $\mathbf{S}^{p-1}$  上的测地距离来衡量两个数据点间的相似程度, 故此时代二元聚类问题可以表示为

$$\arg \min_{r_{nk}, \mathbf{c}_k} J = \sum_{n=1}^N \sum_{k=1}^2 r_{nk} \arccos(\mathbf{x}_n^T \mathbf{c}_k). \quad (8)$$

其中

$$r_{nk} = \begin{cases} 1, & k = \arg \min_j \arccos(\mathbf{x}_n^T \mathbf{c}_j); \\ 0, & \text{otherwise.} \end{cases}$$

$$\mathbf{x}_n, \mathbf{c}_k \in \mathbf{S}^{p-1}.$$

对于优化问题 (8), 可以通过两步迭代的方法进行求解. 首先注意到  $r_{nk}$  是一个二元变量, 通过比较的方法确定, 而  $\mathbf{c}_k$  可以通过连续优化求解. 假设当前迭代中  $r_{nk}$  已确定, 转而考虑优化问题

$$\arg \min_{\mathbf{c}_k} J = \sum_{n=1}^N \sum_{k=1}^2 r_{nk} \arccos(\mathbf{x}_n^T \mathbf{c}_k). \quad (9)$$

由约束优化理论可知, 在最优点处目标函数的常规梯度应平行于数据变量, 即此刻黎曼梯度为零. 计算可得目标函数  $J$  关于  $\mathbf{c}_k$  的常规梯度向量为

$$\partial_{\mathbf{c}_k} J = - \sum_{n=1}^N \frac{r_{nk} \mathbf{x}_n}{\sqrt{1 - (\mathbf{x}_n^T \mathbf{c}_k)^2}}. \quad (10)$$

代入黎曼梯度公式 (6) 可得

$$\nabla_{\mathbf{c}_k} J = - \sum_{n=1}^N \frac{r_{nk} (\mathbf{I}_p - \mathbf{c}_k \mathbf{c}_k^T) \mathbf{x}_n}{\sqrt{1 - (\mathbf{x}_n^T \mathbf{c}_k)^2}}. \quad (11)$$

将式 (11) 置零并整理可得

$$\mathbf{c}_k \sum_{n=1}^N \frac{r_{nk} \mathbf{c}_k^T \mathbf{x}_n}{\sqrt{1 - (\mathbf{x}_n^T \mathbf{c}_k)^2}} = \sum_{n=1}^N \frac{r_{nk} \mathbf{x}_n}{\sqrt{1 - (\mathbf{x}_n^T \mathbf{c}_k)^2}}. \quad (12)$$

以上给出了求取  $\mathbf{c}_k$  的不动点迭代格式. 注意到式 (12) 左端仅是  $\mathbf{c}_k$  与一个比例因子的乘积, 因此在迭代过程中可以仅考虑该比例因子的符号而忽略其数值大小, 最后通过范数归一化达到更新  $\mathbf{c}_k$  的目的. 由式 (12) 也可以发现, 正是黎曼梯度的引入才使得不动点迭代成为可能.

总结以上优化过程可以得到如下单位超球面  $\mathbf{S}^{p-1}$  上的二元聚类算法:

Step 1: 选取初始聚类中心  $\mathbf{c}_k \in \mathbf{S}^{p-1}$ ,  $k = 1, 2$ .

Step 2: 确定每个数据点的类别属性指标

$$r_{nk} = \begin{cases} 1, & k = \arg \min_j \arccos(\mathbf{x}_n^T \mathbf{c}_j); \\ 0, & \text{otherwise.} \end{cases}$$

Step 3: 更新聚类中心

$$\mathbf{c}_k \leftarrow \frac{\sum_{n=1}^N \frac{r_{nk} \mathbf{x}_n}{\sqrt{1 - (\mathbf{x}_n^T \mathbf{c}_k)^2}}}{\text{sign} \left( \sum_{n=1}^N \frac{r_{nk} \mathbf{c}_k^T \mathbf{x}_n}{\sqrt{1 - (\mathbf{x}_n^T \mathbf{c}_k)^2}} \right)}.$$

Step 4: 对  $\mathbf{c}_k$  进行单位化  $\mathbf{c}_k \leftarrow \mathbf{c}_k / \|\mathbf{c}_k\|$ .

Step 5: 如果目标函数  $J$  收敛则停止迭代; 否则返回 Step 2 继续以上过程.

当限定数据类别数  $K = 1$  时, 对于每个数据点  $\mathbf{x}_n$  有  $r_n \equiv 1$ , 这就得到如下求取数据点平均的算法:

Step 1: 选取初始向量  $\mathbf{c} \in \mathbf{S}^{p-1}$ .

Step 2: 更新数据中心

$$\mathbf{c} \leftarrow \frac{\sum_{n=1}^N \frac{\mathbf{x}_n}{\sqrt{1 - (\mathbf{x}_n^T \mathbf{c})^2}}}{\text{sign} \left( \sum_{n=1}^N \frac{\mathbf{c}^T \mathbf{x}_n}{\sqrt{1 - (\mathbf{x}_n^T \mathbf{c})^2}} \right)}.$$

Step 3: 对向量单位化  $\mathbf{c} \leftarrow \mathbf{c} / \|\mathbf{c}\|$ .

Step 4: 如果向量  $\mathbf{c}$  收敛则停止; 否则返回 Step 2 继续以上过程.

下面对上述算法做 4 点说明:

1) 在欧氏空间中实施通常的  $K$ -means 聚类或求取算术平均然后向  $\mathbf{S}^{p-1}$  投影可能被认为是更加直接的方法. 其实不然, 以求取平均为例, 如果数据点在超球面上呈现中心对称分布, 则通常的算术平均将给出退化的解. 此外, 对于更一般的非线性流形凸性条件往往难以得到满足, 因此投影算子不再是适定的.

2) 根据  $\mathbf{S}^{p-1}$  的紧致性, 目标函数  $J$  在  $\mathbf{S}^{p-1}$  上是有界的. 由于在极大值点和极小值点处具有相同的不动点迭代形式, 初始向量的选取将会影响算法实施的结果. 对此可以考虑不动点迭代和梯度下降相结合的学习模式. 在算法实施的初期, 利用梯度下降法将聚类中心移动到极小值点附近; 然后利用不动点迭代求取最终的极小值点. 关于梯度下降法可以利用如下 3 种迭代形式:

$$\mathbf{c}_k \leftarrow \mathbf{c}_k - \lambda \partial_{\mathbf{c}_k} J, \quad \mathbf{c}_k \leftarrow \mathbf{c}_k / \|\mathbf{c}_k\|; \quad (13)$$

$$\mathbf{c}_k \leftarrow \mathbf{c}_k - \lambda \nabla_{\mathbf{c}_k} J, \quad \mathbf{c}_k \leftarrow \mathbf{c}_k / \|\mathbf{c}_k\|; \quad (14)$$

$$\mathbf{c}_k \leftarrow g(\lambda). \quad (15)$$

其中

$$g(t) = \cos(\|\nabla_{\mathbf{c}_k} J\| t) \mathbf{c}_k - \frac{\sin(\|\nabla_{\mathbf{c}_k} J\| t)}{\|\nabla_{\mathbf{c}_k} J\|} \nabla_{\mathbf{c}_k} J$$

是以  $\mathbf{c}_k$  为起点,  $-\nabla_{\mathbf{c}_k} J$  为初始速度的测地线方程,  $\lambda$

是学习步长. 式 (13) 和 (14) 分别是常规梯度流和黎曼梯度流的离散迭代格式, 显然格式 (14) 要优于格式 (13). 而式 (15) 是测地流的离散迭代格式, 测地流技术直接在约束空间中进行优化, 可以省去归一化的步骤, 并具有明显的几何意义.

3) 上述降梯度和不动点相结合的优化策略也保证了算法的局部收敛性, 因为在每次迭代中目标函数的取值是下降的. 文献 [10] 对这种两步迭代算法的收敛性有更详尽的论述. 此外, 由于目标函数是非凸的, 只能给出局部极小值, 通过反复实验或结合一些智能算法可以解决这个问题.

4) 以上的二元聚类算法来源于 ICA 的实际应用背景, 但是可以发现, 上述算法可以直接推广到  $K > 2$  时的数据聚类问题.

## 5 仿真实验与分析

首先考虑超球面二元聚类算法和通常  $K$ -means 聚类算法的性能比较. 在单位超球面上, 数据的概率分布是不可穷尽的, 以下仅考虑最具代表性的 von Mises 分布情形. von Mises 分布是高斯分布在单位超球面投影的结果. 此外, 在  $K$ -means 聚类中, 由于数据的单位范数约束, 还须对其聚类结果进行单位化.

在实验中, 首先随机产生两个原始的中心单位向量; 然后将一组具有相同维数和特定标准差的高斯噪声向量组叠加到中心向量上, 并进行范数归一化. 其中, 两个中心向量所叠加的噪声数据个数也是随机的. 随后分别利用这两种方法进行聚类分析, 且考虑其聚类中心与对应的原始中心向量的相关性差值作为评判这两种方法优劣的标准. 由于有两个聚类中心, 可以考虑相关性差值的平均值, 而相关性差值由球面聚类相关系数减去  $K$ -means 聚类相关系数得到. 表 1 显示了在不同数据维数和噪声标准差下相关性差值的计算结果. 其中, 每次实验的球面样本点均取为 50, 而表中每个数据是重复 100 次实验所得结果的平均值. 从表 1 可以看出, 所有数字都是正值, 这说明球面聚类的效果要优于  $K$ -means 聚类的效果, 且这种优势在低维数据和高水平噪声下越发明显. 尽管表中的每个数据本身是个随机量, 但是如下统计性规律总是存在的, 即除去个别点外, 相关性差值随数据维数的增大和噪声水平的减小而减小.

表 1 相关性差值  $10^{-4}$

| 噪声标准差 | 维数   |      |      |      |      |      |      |
|-------|------|------|------|------|------|------|------|
|       | 5    | 10   | 20   | 30   | 40   | 50   | 60   |
| 0.1   | 2.56 | 2.47 | 1.82 | 1.80 | 1.65 | 1.29 | 1.22 |
| 0.3   | 17.1 | 15.4 | 11.2 | 7.09 | 4.57 | 2.01 | 2.04 |
| 0.5   | 35.5 | 29.0 | 18.7 | 10.2 | 8.12 | 5.98 | 4.23 |

以下考虑该算法在地球物理勘探领域的应用. 依据线性系统假设, 人工地震记录可以看作是地层反射系数序列与地震子波卷积的结果. 由于反射系数与地震子波均未知, 地震反卷积基本上是一个盲过程. 最近, 文献 [11-12] 提出了一种称为带状独立成分分析 (BICA) 反卷积的方法. 该方法将单道反卷积问题转化为多通道 ICA 问题; 然后利用 ICA 中的算法对反卷积问题实施求解. 相对于经典的最小熵和 Bussgang 型反卷积方法, BICA 反卷积收敛更快并具有较强的鲁棒性.

利用 ICA 的并行处理结构可以分解出与子波长度数值相等的多道候选子波, 为此, 文献 [11] 提出了一种基于最小二乘的选取方法, 但是这种方法对于子波的时移相当敏感. 此外, 在人工地震勘探中一般认为地震子波作为系统输入在一定的空间范围内具有相同或相近的形式, 因此可以利用邻近道间的相关信息来克服子波提取中时移的影响. 为了使学习过程稳定, 假设所有子波均具有归一化能量, 这样就自然赋予了子波单位超球面的几何结构. 由于 ICA 符号的不确定性, 提取出的子波向量在超球面上呈现出近似对称的二元分布结构. 可以考虑在单位超球面上对多道子波实施二元聚类, 然后将其中一类数据和另一类数据的反号合并, 并求取这一类新数据的球面平均, 以此消除测量系统和单道数据处理中的随机影响. 子波提取不是本文关注的重点, 相关细节见文献 [11-12].

以下采用采样数为 400 点且满足拉普拉斯分布的序列模拟地层反射系数, 而地震子波采用零相位雷克子波模拟<sup>[11]</sup>, 且采样数为 36 点, 两者卷积后就得到模拟地震数据. 经过 BICA, 可以对提取出的多道子波实施球面二元聚类和求取平均. 图 2 显示了对 25 道模拟地震数据实验的结果. 其中, 图 1(a) 是卷积用的 36 点零相位雷克子波, 图 1(b) 是对基于邻近道相关所提取出的 25 道子波实施超球面二元聚类和平均的结果, 图 1(c) 是对这 25 道子波实施欧氏聚类和平均的结果, 图 1(d) 是对基于文献 [11] 中最小二乘法所提取出的 25 道子波实施欧氏聚类和平均的结果. 由图 2 可见, 利用单位超球面上聚类和平面的方法提取出的地震子波明显优于直接用欧氏聚类和算术平均所得的结果, 且图 1(c) 中的子波在时间轴上被拉伸, 因而主频降低, 利用其实施反卷积势必将降低信号的分辨率, 图 1(d) 的结果显然是不能被接受的.

改变地震数据的道数, 进行子波提取可以得到相似的结果. 图 3 显示了道数分别为 5, 10, 15, 20, 25 和 30 时, 用这两种方法提取出的子波与原始雷克子波间相关系数的比较, 黑色标志为超球面聚类及平均, 灰色标志为欧氏聚类及平均. 由于时移影响, 这些相

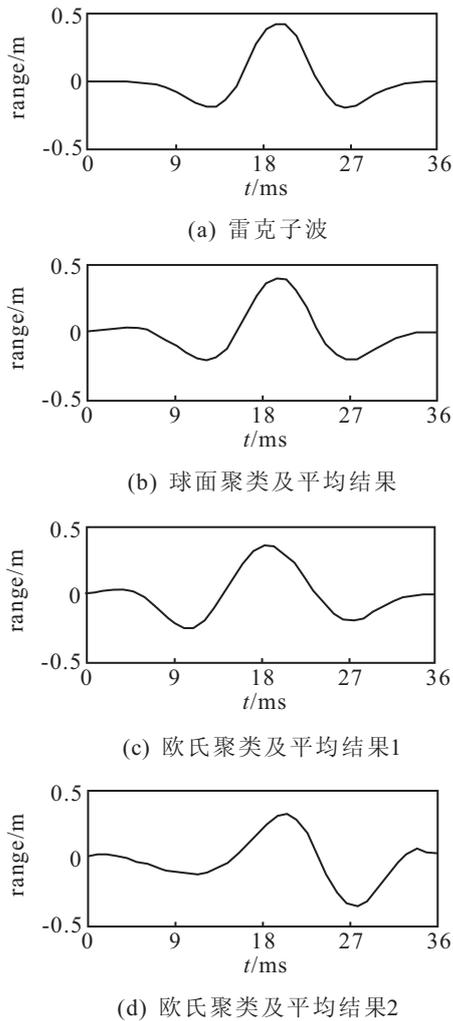


图 2 各种子波的比较

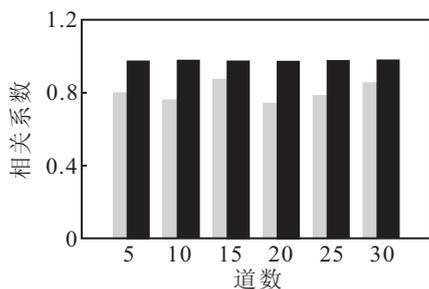


图 3 相关系数的比较

关系数是依据子波振幅谱计算而得的,且图中的数值结果是进行 50 次实验的平均值.由图 3 可以看出,基于超球面聚类和平均的子波提取方法具有明显的优势,在提取过程中几乎完全恢复了子波的振幅谱信息.

## 6 结 论

本文利用黎曼度量研究了单位超球面上二元聚类和求取数据平均的方法,然后将其应用于地震子波提取并取得了较通常方法更为满意的结果.而且,这种方法能够直接推广到单位超球面上任何聚类分析

或数据学习算法上.本文只研究了 von Mises 分布下这两种聚类分析的比较,对于一般情况下的性能比较还值得进一步研究.此外,由于单位超球面是一个常曲率空间,它同欧氏空间的特征维数仅相差 1,其弯曲程度并不是很明显,而在欧氏空间中更低维嵌入子流形上的机器学习理论是一个值得关注和研究的方面.

## 参考文献(References)

- [1] Roweis S T, Saul L K. Nonlinear dimensionality reduction by locally linear embedding[J]. Science, 2000, 290(22): 2323-2326.
- [2] Tenenbaum J B, Silva de V, Langford J C. A global geometric framework for nonlinear dimensionality reduction[J]. Science, 2000, 290(22): 2319-2323.
- [3] Amari S, Fiori S. Geometrical methods in neural networks and learning[J]. Neurocomputing, 2005, 67(8): 1-7.
- [4] Goh A, Vidal R. Unsupervised Riemannian clustering of probability density functions[C]. European Conf on Machine Learning and Principles and Practice of Knowledge Discovery in Databases. Antwerp, 2008: 377-392.
- [5] Goh A, Vidal R. Clustering and dimensionality reduction on riemannian manifolds[C]. IEEE Conf on Computer Vision and Pattern Recognition. Anchorage, 2008: 1-7.
- [6] Kim J, Shim K H, Choi S. Soft geodesic kernel  $K$ -means[C]. IEEE Int Conf on Acoustics, Speech and Signal Processing. Honolulu, 2007, 2: 429-432.
- [7] Hyvarinen A, Karhunen J, Oja E. Independent component analysis[M]. New York: Wiley, 2001.
- [8] Wang F S, Li H W, Li R. Data mining with independent component analysis[C]. The 6th World Congress on Intelligent Control and Automation. Dalian, 2006, 2: 6043-6047.
- [9] Neill B O. Elementary differential geometry[M]. New York: Academic, 1966.
- [10] MacQueen J B. Some methods for classification and analysis of multivariate observations[C]. Proc of 5th Berkeley Symposium on Mathematical Statistics and Probability. Berkeley, 1967, 1: 281-297.
- [11] Kaplan S T, Ulrych T J. Blind deconvolution and ICA with a banded mixing matrix[C]. The 4th Int Symposium on ICA and Blind Signal Separation. Nara, 2003: 223-228.
- [12] Aws A Q, Woo W L, Dlay S S. Blind seismic deconvolution of single channel using instantaneous independent component analysis[C]. The 6th Int Symposium on Communication Systems Networks and Digital Signal Processing. Graz, 2008: 142-146.