

黎曼法坐标流形学习扩展算法*)

韦佳 彭宏 林毅申

(华南理工大学计算机科学与工程学院 广州 510640)

摘要 LOGMAP 是最近提出的一种黎曼流形学习算法,它能够有效地学习出高维数据的低维嵌入坐标。然而该算法只能处理单类数据的流形学习问题,当存在多类数据时往往不能得到理想的嵌入结果。为解决这个问题,提出了一种扩展的 LOGMAP 算法(Extended LOGMA PAlgorithm,简称 ELOGMAP)。该算法通过计算全局基准点所在类到其他类的最短距离找出各类的局部基准点,然后逐个计算各类数据相对于局部基准点的局部黎曼法坐标,最后通过扩展的全局基准点与局部基准点之间测地距离关系得到多类数据的整体嵌入坐标。实验结果验证了该算法在处理多类数据流形学习上的有效性。

关键词 流形学习,黎曼法坐标,对数映射

Extended Manifold Learning Algorithm Based on Riemannian Normal Coordinates

WEI Jia PENG Hong LIN Yi-shen

(School of Computer Science and Engineering, South China University of Technology, Guangzhou 510640, China)

Abstract LOGMAP is a Riemannian manifold learning algorithm proposed recently. It is efficient for many nonlinear dimension reduction problems. However, the algorithm is only fit for single class problem, when applied to multi-class data, it can't get desirable embedding. In this paper, an extended LOGMAP algorithm is proposed for the multi-class manifold learning. In the algorithm, the local base points are found through the shortest distance from the class of global base point to other classes. Then the local Riemannian normal coordinates of each class are calculated based on the local base point of the class. At last, the global low-dimensional coordinates are obtained by the relationship of extended geodesic distance between the global base point and the local base points. Experimental results demonstrate the validity of the algorithm.

Keywords Manifold learning, Riemannian normal coordinates, Logarithmic map

当前计算机技术的许多领域,如图像分析与计算机视觉、生物信息学和生物特征识别、信息检索及文本挖掘等,其数据都具有高维数的特点。数据的这种特点一方面会导致维数灾难的出现,使得已有的方法难以对其进行有效处理;另一方面也使得人们难以感知和理解这些数据,不能充分发现和利用它们的内在结构。因此,必须对高维数据进行维数约减,这是模式识别、机器学习、数据挖掘等领域的重要研究课题之一^[1]。

传统的维数约减方法如主成分分析(PCA, Principal Component Analysis)与多维尺度变换(MDS, Multidimensional Scaling)能够有效地学习出具有线性结构的高维数据的内在结构。但是,实际应用中很多数据在高维空间中常常是高度非线性、属性强相关的,在这种情况下,传统的维数约减方法显然是无能为力的。

流形学习为解决上述问题提供了一条可行的思路。流形可以看成是高维空间中的曲面,从局部来看它与欧式空间是同胚的,但是从整体来看它却是弯曲的并具有全局非线性的特点。在任一点都定义有黎曼度量张量的微分流形称之为黎曼流形^[2]。流形学习的目的是当数据集所在的空间表现为流形结构时,要从这些数据集中恢复出流形的内在几何结构及其规律性^[3]。它是从观测到的现象中寻找事物的本质,找到数据产生的内在规律。要达到这一目标,必须对低维嵌入做

适当的约束或限制,不同的约束和限制也就得到不同的流形学习算法^[4]。目前具有代表性的流形学习算法有 LLE^[5], ISOMAP^[6], LTSA^[7], SDE^[8], LOGMAP^[9]等。

上面所列的流形学习算法都能够有效地学习出体现数据集内在低维流形结构的整体嵌入坐标,而且各有各的特色与优点。但是这些流形学习算法都面临着同样的一个问题,这就是这些算法只是考虑了在只有一类数据的情况下如何发现数据集的内在低维流形结构,而没有考虑多类数据同时存在的情况下如何发现有效的低维流形嵌入。SLLE^[10]和 EISOMAP^[11]为解决这个问题做了一些有益的尝试,但是这些方法更多地是注意保持了类间数据的相互关系,没有能够很好地恢复同类数据之间的内在流形结构。基于此,本文提出了一种扩展的 LOGMAP 算法 ELOGMAP(Extended Logarithmic Ma PAlgorithm)来处理多类数据情况下的低维流形嵌入问题。该算法既能保持类间数据的区别,又能够恢复类内数据的流形结构。在模拟数据集和真实数据集上的实验结果表明,ELOGMAP 算法能较好地处理多类数据情况下的流形学习问题。

1 LOGMAP 算法简介

1.1 黎曼法坐标

*)广东省自然科学基金项目(07006474),广东省科技攻关项目(2007B010200044)。韦佳 博士研究生,研究方向为人工智能、机器学习;彭宏 教授,博导,研究方向为神经网络、遗传算法、数据挖掘;林毅申 博士研究生,研究方向为进化算法、基因表达式编程。

流形上某一点的黎曼法坐标包含了该点到流形上一特定点的距离与方向的信息。黎曼法坐标的几何直观解释可以通过图1来说明^[9]。假设M代表的是地球的表面， p 是地球上的一特定点，那么朝着某一方向并沿着测地线走，我们可以到达地球上的任何一个地方。如果图中的 a, b, c 三点代表地球上三个不同的地方，那么根据它们距 p 点的距离与方向的信息，可以把三维空间中球面上的点映射到二维平面上，也就是图右侧的 a', b', c' ，这就是 a, b, c 的黎曼法坐标。该映射为指数映射(Exponential Map)的逆映射，称之为对数映射(LOG-MAP, Logarithmic Map)。

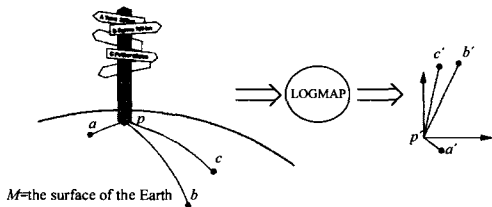


图1 黎曼法坐标示意图

1.2 LOGMAP 算法介绍

LOGMAP算法的基本思想是用低维空间的黎曼法坐标来表示高维流形上的点，实现维数约减的目的^[9]。即通过对数映射发现高维数据集 $X = \{x_1, x_2, \dots, x_N\}, x_i \in R^D$ 在低维空间中的坐标表示 $Y = \{y_1, y_2, \dots, y_N\}, y_i \in R^d (d < D)$ 。该算法主要包括四步，描述如下：

输入：样本集 X ，基准点 p ，低维嵌入空间的维数 d ，样本点的近邻数 K ；

输出： X 的低维嵌入坐标 Y 。

第一步 根据基准点 p ，计算其 K 最近邻，记为 $N(p)$ 。

第二步 计算 $N(p)$ 的低维坐标及 p 点处切空间的标准正交基。由于这些点是近似线性的，可以通过标准的PCA或MDS算法得到它们在 R^d 空间中的低维坐标及 p 点处切空间的标准正交基。

第三步 计算 $N(p)$ 中的点到其他各点的测地距离。首先在高维空间构造近邻图，寻找和每个样本距离最近的 K 个样本作为近邻，每个样本和它的近邻之间建立连接，从而构成一个无向图。然后计算 $N(p)$ 中的点到其他各点的测地距离，这一步可以用Dijkstra算法求解，利用两点之间最短路径近似测地距离。

第四步 根据上面所得信息，计算 $N(p)$ 之外其它各点的黎曼法坐标。

算法中第四步的黎曼法坐标求解有很多种方法^[12]，这里介绍一种通过二阶多项式插值求解黎曼法坐标的方法。假设流形上任一点所处的区域都是凸区域，那么对数映射可用下式表示^[13]：

$$\text{Log}_p(x_i) = -\frac{1}{2} \nabla_x d^2(x, x_i) |_{x=p} \quad (1)$$

又假设 $N(p)$ 外一点 x' 到 $N(p)$ 中各点 x 的距离的平方 $d^2(x, x')$ 可由二阶多项式表示：

$$h(y) = a + b_i y^i + C_{*} y^i y^i \quad (2)$$

式中 $y \in R^d$ 为 $x \in R^D$ 的低维坐标， y^i, y^j 为 y 的分量， $b_i y^i$ 和 $C_{*} y^i y^i$ 为符合爱因斯坦求和约定的多项式。上式对某一分量 y^k 的偏导数为

$$\frac{\partial h(y)}{\partial y^k} = b_k + C_{*k} y^i + C_{*k} y^i \quad (3)$$

由上述描述可知，问题归结为求解出 a, b, C_{*} ，使得下式成立。

$$\begin{pmatrix} 1 & y_1^i & y_1^i y_1^i \\ 1 & y_2^i & y_2^i y_2^i \\ \dots & \dots & \dots \\ 1 & y_k^i & y_k^i y_k^i \end{pmatrix} \begin{pmatrix} a \\ b_i \\ C_{*i} \end{pmatrix} = \begin{pmatrix} d^2(x_1, x') \\ d^2(x_2, x') \\ \dots \\ d^2(x_k, x') \end{pmatrix} \quad (4)$$

由此可得

$$\begin{pmatrix} a \\ b_i \\ C_{*i} \end{pmatrix} = \begin{pmatrix} 1 & y_1^i & y_1^i y_1^i \\ 1 & y_2^i & y_2^i y_2^i \\ \dots & \dots & \dots \\ 1 & y_k^i & y_k^i y_k^i \end{pmatrix}^{-1} \begin{pmatrix} d^2(x_1, x') \\ d^2(x_2, x') \\ \dots \\ d^2(x_k, x') \end{pmatrix} \quad (5)$$

其中右边第一项为摩尔-彭诺斯逆，即加号逆。由求解出的 a, b_i 与 C_{*} ，根据偏导数公式(3)及对数映射公式(1)就能求出 x' 的黎曼法坐标。图2所示为LOGMAP算法对S-Curve曲面的降维结果，从图中可以看出LOGMAP算法正确地发现了S-Curve曲面的内在流形结构。

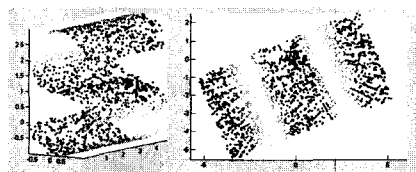


图2 S-Curve曲面及其LOGMAP降维结果(黑色点表示 $N(p)$)

2 ELOGMAP 算法

2.1 问题描述

在LOGMAP算法中，第三步是至关重要的，它的基本思想是用两点之间的最短距离来近似这两点之间的测地距离。然而，如果存在多类数据，情况就不同了。以图3(a)为例，图中两个不连续的曲面分别代表两类数据。如果把基准点选在图中上半部分所示的数据类中(记为 C_1 类，黑色点代表基准点及其近邻 $N(p)$)，那么黑色点到 C_1 类的其它点的测地距离可以很容易求到。然而，黑色点到下半部分所示的数据类(记为 C_2 类)的测地距离都为无穷大，这样的话就不能求出下半部分数据点的黎曼法坐标嵌入。

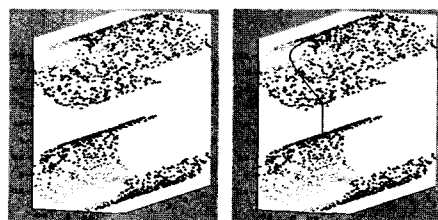


图3 两个半S-Curve数据集

2.2 ELOGMAP 算法介绍

为了解决上面所提出的问题，本节提出一种扩展的LOGMAP算法(ELOGMAP, Extended Logarithmic Map Algorithm)来解决多类数据的流形学习问题。假设一共有 N 类数据，算法描述如下：

输入：样本集 X ，全局基准点 p ，低维嵌入空间的维数 d ，样本点的近邻数 K ，样本的类别标记 ClassIndex ，样本的类别数 N ；

输出： X 的低维嵌入坐标 Y 。

第一步 根据全局基准点 p , 计算其 K 近邻点 $N(p)$ 。假设 p 点属于 C_1 类, 那么只计算与 p 点同类的点的黎曼法坐标, 这一步与原始 LOGMAP 算法一样。

第二步 找出两个分别属于 C_1 类与 $C_i (i \neq 1)$ 类数据的点, 使得这两点之间的欧式距离是这两类数据之间最短的。令该距离为 $d_{c_1 c_i}$, 即找出点 $i \in C_1$ 与点 $j \in C_i$, 使得 $d_{c_1 c_i} = \min_{i \in C_1, j \in C_i} d\{i, j\}$ 。那么 $N(p)$ 中各点到 j 点的测地距离可由下式得到:

$$d_G(N(p), j) = d_G(N(p), i) + d(i, j) \quad (6)$$

如图 3(b) 所示, 蓝色曲线代表 $d_G(N(p), i)$, 绿色直线代表 $d(i, j) = d_{c_1 c_i}$ 。

第三步 根据对数映射公式, 计算 j 点相对于 p 点的黎曼法坐标。

第四步 以 j 点为局部基准点, 计算与 j 点同类的的数据即 C_i 类数据相对于 j 的黎曼法坐标。根据 j 点相对于 p 点的坐标得到 C_i 类的全局嵌入坐标。

第五步 重复第二步至第四步, 直到所有 N 类数据黎曼法坐标计算完为止。

需要注意的是, 因为该扩展算法是以选定的基准点 p 为原点计算各类数据的黎曼法坐标, 所以非基准类的数据的低维坐标可能会相互重叠, 影响算法功能及可视化效果。为了防止这种情况发生, 需要修改 $d_{c_1 c_i}$ 的计算公式。本文采用如下方法来解决这个问题: 把已经得到低维嵌入坐标的数据看成一个整体, $maxlength_p$ 代表 p 点到这个整体中的数据点的最大测地距离; 把将要计算低维嵌入坐标的数据类 C_i 看成另外一个整体, $maxlength_i$ 代表 C_i 类中离基准类 C_1 最近的 j 点到 C_i 类其它点的最大测地距离。那么令

$$d_{c_1 c_i} = \min_{i \in C_1, j \in C_i} d\{i, j\} + maxlength_p + maxlength_i \quad (7)$$

用这个公式代替第二步中的原公式, 就可以避免非基准类的数据的低维坐标相互重叠的情况。

3 实验及分析

为了验证算法的有效性, 本节给出了 ELOGMAP 算法在模拟数据集及 MNIST 数据集上的实验结果及分析。

3.1 模拟数据

本实验使用 2.1 节所示的模拟数据集, 该数据集由两个半 S-Curve 数据集组成, 随机均匀采集 2000 个数据点, 每个半 S-Curve 中有 1000 个数据点, 这样该问题就可以看成是两类数据的流形学习问题。在实验中, 我们取 $d=2, K=30$ 。图 4 展示了基准点 p (黑色点表示 $N(p)$) 取两个不同位置时的实验结果。从图中可以看出, ELOGMAP 算法很好地解决了两类数据的流形学习问题, 不仅学习到了每类数据的内在流形结构, 而且保持了不同类数据之间的区别, 取得了较好的可视化效果。



(a) p 点在 C_1 类中 (b) p 点在 C_2 类中
图 4 两个半 S-Curve 在二维空间的降维结果

3.2 MNIST 数据集

MNIST 数据集 (<http://yann.lecun.com/exdb/mnist>) 由 10 个手写数字的 60000 个训练样本和 10000 个测试样本组成, 每个样本都是一幅 28×28 的灰度图像, 如图 5 所示。本实验从 60000 个训练样本中随机选取数字 '0', '1', '2' 的图像各 1000 幅, 对它们用 ELOGMAP 算法降维, 该问题可以看成是三类数据集的流形学习问题。实验中, 取 $d=2, K=40$ 。图 6(a) 是 ELOGMAP 算法只对数字 '1' 降维的结果, 图 6(b) 是 ELOGMAP 算法对三类数字降维的实验结果。从图 6(b) 中可以看出, 这三类数字的低维嵌入明显聚集为 3 类, 而且每类数字在低维坐标中的分布都有着极强的规律性。以数字 '1' 为例, 随着横坐标的增加, 数字的倾斜角度由向左倾斜渐变为向右倾斜; 随着纵坐标的增加, 数字的宽度由粗变细。这与图 6(a) 所示数字 '1' 的排列规律相似, 只是坐标轴的方向不同而已。可见 ELOGMAP 算法的确能学习到多类 MNIST 数据集的潜在低维变化规律。

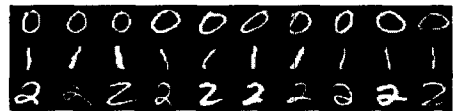
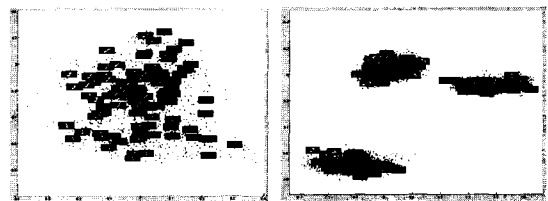


图 5 MNIST 数据集



(a) 单类数字降维结果 (b) 三类数字降维结果

图 6 三类手写数字由 ELOGMAP 计算所得的二维坐标

结束语 本文提出了一种 LOGMAP 的扩展算法 ELOGMAP, 该算法克服了原始 LOGMAP 算法只能处理一类数据的流形学习的缺陷, 可以处理多类数据的流形学习问题。ELOGMAP 算法采用“芝麻开花”的形式逐个计算各类数据的黎曼法坐标, 既保持了类内数据的流形结构, 又保持了类间数据的相互关系, 使得低维坐标不致于混合重叠。实验结果表明, ELOGMAP 算法能有效地获取多类数据的低维嵌入, 取得了较好的可视化效果。

ELOGMAP 算法也存在一些需要改进和进一步研究的问题。比如说, ELOGMAP 仅仅是简单地通过距离来保持类间的关系, 是否存在更好的方法来保持这种关系; 另外, 当数据维数较高, 类别较多的时候, 算法运行速度较慢, 寻找更快速简单的方法, 以及将该算法用于模式识别领域等, 都是我们需要努力的方向。

参考文献

- [1] 罗四维, 赵连伟. 基于谱图理论的流形学习算法. 计算机研究与发展, 2006, 43(7): 1173-1179
- [2] 陈维恒. 微分流形初步(第二版). 北京: 高等教育出版社, 2001
- [3] 张军平. 流形学习及其应用. 北京: 中国科学院研究生院, 2003
- [4] 杨剑. 流形学习若干问题研究. 北京: 中国科学院研究生院, 2006
- [5] Roweis S T, Saul L K. Nonlinear Dimensionality Reduction by Locally Linear Embedding. Science, 2000, 290(5500): 2323-2327

(下转第 209 页)

- (4) $\bar{S} \oplus (F \vee G) \geq \bar{S} \oplus F \vee \bar{S} \oplus G$;
- (5) $S \otimes (F \vee G) = S \otimes F \vee S \otimes G$;
- (6) $S \otimes (F \wedge G) \leq S \otimes F \wedge S \otimes G$;
- (7) $\bar{S} \oplus U = U$ 并且 $S \otimes \Phi = \Phi$.

定理 4.10 设 U 是论域, $F: U \rightarrow [0, 1]$ 是 U 上的模糊集, $S, G: U \times U \rightarrow [0, 1]$ 是 U 上的模糊关系, F, S, G 分别是 F, S, G 的布尔矩阵, 则有

- (1) $\overline{(S \vee G)} \oplus F = \bar{S} \oplus F \wedge \bar{G} \oplus F$;
- (2) $\overline{(S \vee G)} \otimes F = S \otimes F \vee G \otimes F$;
- (3) $\overline{(S \otimes G)} \oplus F = \bar{S} \oplus (\bar{G} \oplus F)$;
- (4) $\overline{(S \otimes G)} \otimes F = S \otimes (G \otimes F)$;
- (5) $S^* \oplus F = F \wedge (\bar{S} \oplus (S^* \oplus F))$;
- (6) $S^* \otimes F = F \vee (S \otimes (S^* \otimes F))$.

定理 4.11 设 $U = (e_1, \dots, e_n)$ 是论域, $F: U \rightarrow [0, 1]$ 是 U 上的模糊集, $S: U \times U \rightarrow [0, 1]$ 是 U 上的模糊等价关系, $F = (F(e_1), \dots, F(e_n))$, $S = (s_{ij})$ 分别是 F, S 的布尔矩阵, 则有

- (1) $\bar{S} \oplus F \leq F \leq S \otimes F$;
- (2) $\bar{S} \oplus F = \bar{S} \oplus (\bar{S} \oplus F)$;
- (3) $S \otimes F = S \otimes (S \otimes F)$;
- (4) 如果 $\forall i \forall j (s_{ij} \vee \bar{s}_{ij} \geq F(e_i))$, 则有 $S \otimes F = \bar{S} \oplus (S \otimes F)$;
- (5) 如果 $\forall i \forall j (s_{ij} \wedge \bar{s}_{ij} \leq F(e_i))$, 则有 $\bar{S} \oplus F = S \otimes (\bar{s} \oplus F)$.

结束语 本文是矩阵论方法在模糊粗糙集理论研究中的具体应用, 其主要目的是运用矩阵这一有力的数学工具, 对模糊粗糙集的基本概念和基本运算性质给出一种较为系统和完整的描述。在布尔矩阵逻辑运算中, 同时定义“与积”和“或积”两种运算, 较好地实现了上述目的, 尤其是针对模糊粗糙集理论中的模糊必然算子和模糊可能算子计算过程的布尔矩阵表示, 为基于模糊粗糙集理论的知识表示与知识获取提供了一种能行与可计算的思路与方法。

致谢 感谢高尚博士为本文提供资料和在成文过程中的积极建议。

参 考 文 献

[1] Pawlak Z. Rough Sets[J]. International Journal of Computer

(上接第 200 页)

[6] Tenenbaum J B, de Silva V, Langford J C. A Global Geometric Framework for Nonlinear Dimensionality Reduction. Science, 2000, 290(5500): 2319-2323

[7] Zhang Zhenyue, Zha Hongyuan. Principal Manifolds and Nonlinear Dimensionality Reduction via Tangent Space Alignment. SIAM Journal of Scientific Computing, 2004, 26(1): 313-338

[8] Weinberger K, Saul L. Unsupervised Learning of Image Manifolds by Semidefinite Programming // Proceedings of IEEE Conference on Computer Vision and Pattern Recognition. Washington, DC, 2004

[9] Brun A, Westin C-F, Herberthson M, et al. Fast Manifold Learning Based on Riemannian Normal Coordinates // Proceedings of the 14th Scandinavian Conference on Image Analysis. Joensuu, Finland, 2005

and Information Science, 1982, 11(5): 341-356

[2] 王国胤. Rough 集理论与知识获取[M]. 西安: 西安交通大学出版社, 2001

[3] 曾黄麟. 粗集理论及其应用—关于数据推理的新方法. 修订版[M]. 重庆: 重庆大学出版社, 1998

[4] Dubois D, Prade H. Rough fuzzy sets and fuzzy rough sets [J]. International Journal of General Systems, 1990, 17: 191-209

[5] Greco S, Matarazzo B, Slowinski R. Fuzzy similarity relation as a basis for rough approximations [A] // The Proceedings of RSCTC. Heidelberg: Springer-Verlag, 1998: 283-289

[6] Morsi N N, Yakout M M. Axiomatics for fuzzy roughsets [J]. Fuzzy Sets and Systems, 1998, 100: 327-342

[7] Radzikowska A M, Kerre E E. A comparative study of fuzzy rough sets [J]. Fuzzy Sets and Systems, 2002, 126: 137-155

[8] Wu W Z, Mi J S, Zhang W X. Generalized fuzzy rough sets [J]. Information Sciences, 2003, 151: 263-282

[9] Mi J S, Zhang W X. An axiomatic characterization of a fuzzy generalization of rough sets [J]. Information Sciences, 2004, 160: 235-249

[10] 黄正华, 胡宝清. 模糊粗糙集理论研究进展 [J]. 模糊系统与数学, 2005, 19(4): 125-134

[11] 苗夺谦, 王珏. 粗糙集理论中概念与运算的信息表示 [J]. 软件学报, 1999, 10(2): 113-116

[12] 李龙星, 运士伟, 杨炳儒. 粗糙集概念与运算的布尔矩阵表示 [J]. 计算机工程, 2005, 31(14): 16-17

[13] 仁艳玲, 朱明放. 基于粗糙集的属性约简的矩阵方法 [J]. 陕西理工学院学报, 2006, 22(3): 76-80

[14] 张桂芸, 黄国兴, 杨炳儒. 基于分辨相似矩阵的相似粗糙集的属性约简算法 [J]. 计算机工程, 2006, 32(10): 43-44

[15] 高学军, 丁军. 基于简化差别矩阵的属性约简算法 [J]. 系统工程理论与实践, 2006, 20(6): 101-107

[16] 雷晓蔚. 粗集理论的矩阵方法 [J]. 计算机工程与应用, 2006, 42(17): 73-75

[17] 张晓如, 张再跃. 基于特征矩阵的粗糙集代数运算与表示定理 [J]. 计算机科学, 2008, 35(4): 170-173

[18] Thiele H. Fuzzy rough sets versus rough fuzzy sets—an interpretation and a comparative study [R]. Technical Report CI-30/98. University of Dortmund, 1998

[10] de Ridder D, Kouropteva O, Okun O, et al. Supervised Locally Linear Embedding // Proceedings of ICANN/ICONIP 2003. LNCS, 2003, 2714: 333-341

[11] Wu Yiming, Chan K L. An Extended Isomap Algorithm for Learning Multi-class Manifold // Proceedings of IEEE International Conference on Machine Learning and Cybernetics (ICMLC2004). Shanghai, China, 2004

[12] Lin Tong, Zha Hongbin, Lee S U. Riemannian Manifold Learning for Nonlinear Dimensionality Reduction // Proceedings of the 9th European Conference on Computer Vision. Graz, Austria, 2006

[13] Fletcher P T, Lu C, Pizer S M, et al. Principal Geodesic Analysis for the Study of Nonlinear Statistics of Shape. IEEE Transactions on Medical Imaging, 2004, 23(8): 995-1005