

基于有监督流形学习的正交投影降维

蒋 润, 周激流, 雷 刚, 李晓华

(四川大学计算机学院, 成都 610064)

摘要: 将监督局部线性嵌入的思想引入传统的正交投影降维方法(OPRA)方法, 提出一种新的基于有监督流形学习的正交投影降维方法(α -OPRA), 使高维到低维的映射在保留某些流形结构的同时, 进一步获得较好的正交投影效果。该方法通过加入额外的参数 α 来控制监督的程度, 在纯粹的有监督的 OPRA 和无监督的 OPRA 之间取得了某些折中。实验结果证明, 该方法能获得较好的降维结果。

关键词: 正交投影降维方法; 降维; 人脸识别

α -based Supervised Orthogonal Projection Reduction by Affinity

JIANG Run, ZHOU Ji-liu, LEI Gang, LI Xiao-hua

(School of Computer, Sichuan University, Chengdu 610064)

【Abstract】 This paper introduces the idea of SLLE into the traditional method of OPRA, which proposes a new approach of α -based Supervised Orthogonal Projection Reduction by Affinity(α -OPRA) for dimension reduction. Such method keeps the reservations of some flow-shaped structure during high-dimensional to low-dimensional mapping, gets better orthogonal projection. The method by adding additional parameters to control the degree of supervision, so in a purely supervised OPRA and unsupervised OPRA between there has been some compromise. Experimental results show that this method can get better reduction result.

【Key words】 Orthogonal Projection Reduction by Affinity(OPRA); dimension reduction; face recognition

1 概述

降维算法是目前特征提取、模式分析、数据挖掘中最为强有力的工具, 对降维算法的研究具有很高的学术价值和实用潜力。其中, 通过降维方法获取人脸子空间, 是目前人脸识别领域最为常用的一种特征提取方法。目前人脸识别领域主流的降维方法有: (1)基于特征向量的方法, 如主成分分析^[1](Principal Component Analysis, PCA); (2)基于流形学习的方法, 如局部线性嵌入^[2](Locally Linear Embedding, LLE); (3)基于变换的方法; (4)基于核的特征提取方法; (5)基于模型的方法。

PCA 是对空间中可变因素的线性关系进行建模, 从而实现高维到低维的映射。传统的 PCA 方法只考虑了数据的二阶统计信息, 未能利用数据中的高阶统计信息, 忽略了空间中数据的非线性相关性。因此, 利用 PCA 进行降维, 将会导致空间中数据信息的丢失。LLE 是基于几何直觉的无监督流形学习的降维方法, 目前在人脸表情识别方面运用较多。在传统的 LLE 方法中, 通过度量欧氏距离的方法可找到每个数据点的 K 个最近邻数据点。利用每个数据点的 K 个最近邻数据点对原始数据进行重构。但是, 在典型流形学习的降维方法中, 缺少对外来样本向低维嵌入空间的映射方法, 这是目前流形学习降维方法普遍存在的问题。

针对 PCA 和 LLE 的不足, 文献[3]通过将 LLE 和 PCA 的思想融合, 提出一种基于流形学习的正交投影降维方法(Orthogonal Projection Reduction by Affinity, OPRA), 并通过理论推导和实验数据证明了其有效性; 文献[4]将 OPRA 降维方法与拉普拉斯特征提取方法相结合进行人脸识别, 并取得了较高识别率。但传统的 OPRA 是无监督的, 没有利用到各数据点的类别信息。虽然文献[3]结尾提到了一种有监督的

OPRA, 但该方法规定每个数据点的 K 个最近邻点必须为相同类别的点, 较为原始, 容易导致空间流形结构的丧失。本文在传统 OPRA 的基础上, 提出了一种新的基于有监督流形学习的正交投影降维算法(α -Based Supervised Orthogonal Projection Reduction by Affinity, α -OPRA)。该方法通过加入额外的参数来控制监督的程度, 在纯粹的有监督的 OPRA 和无监督的 OPRA 之间取得了某些折中, 使高维到低维的映射在保留某些流形结构的同时, 也进一步获得了更好的正交投影效果。

2 传统OPRA降维方法

传统的 OPRA 降维方法引入原始流形学习降维算法——局部线性嵌入, 兼具有 PCA 和 LLE 算法的优点, 并相互弥补了其对方在特征提取时的不足。传统的 OPRA 算法描述为:

(1)利用 PCA 方法求解特征空间矩阵 $V_{PCA} \in R^{m \times (n-c)}$, 将 X 降至 $n-c$ 维, 其中, X 为样本集, $X \in R^{m \times n}$; m 为样本维数; n 为样本总数; c 为样本类别数。

(2)计算每个数据点 $x_i, i=1, 2, \dots, n$ 的 K 个最近邻数据点 $\eta_{ik}, k=1, 2, \dots, K$ 。

(3)利用最近邻数据点计算最优重构权重系数矩阵 W , $W = \{\omega_j\}, j=1, 2, \dots, n$ 。

$$\omega_j = \frac{\sum_k C_{jk}^{-1}}{\sum_{lm} C_{lm}^{-1}} \quad (1)$$

其中, $C_{lm} = (x_j - \eta_{jl}) \cdot (x_j - \eta_{jm}), C \in R^{K \times K}$, 权值 W_{ij} 说明第 j 个数据点对重构第 i 个数据点的所做的贡献^[2]。

作者简介: 蒋 润(1984-), 男, 硕士研究生, 主研方向: 模式识别; 周激流, 教授、博士生导师; 雷 刚, 博士研究生; 李晓华, 副教授
收稿日期: 2009-06-10 **E-mail:** lxhw@scu.edu.cn

(4) 计算 $M = (I - W^T)(I - W)$, $\tilde{M} = XMX^T$ 。

(5) 计算 $y_i = V^T x_i$, 其中, $V = V_{PCA} V_{OPRA}$, $i = 1, 2, \dots, n$, $V \in R^{m \times d}$, V_{OPRA} 为 \tilde{M} 矩阵最小特征值对应特征向量的前 $2:d+1$ 维(舍弃最小特征值对应特征向量), d 为降维后的维数。

3 α -OPRA降维方法

传统的 OPRA 融入了 LLE 的思想, 而 LLE 属于原始的流形学习降维算法。原始的流形学习降维算法都属于无监督的方法, 追求的是挖掘数据本身的内部结构, 通常用于数据挖掘和数据的可视化, 在这些过程中数据的类别信息以及各类之间的关系是未知的, 并没有利用到各数据点的类别信息。文献[3]结尾提到了一种简单的有监督 OPRA 方法(S-OPRA), 规定每个数据点的邻接点仅从与其类别相同的点集中选取, 因此, 邻接点的寻找过程中已经包含了分类的先验知识, 这样有监督的方法更适合于模式识别问题。但实验表明, 这种监督学习较为原始, 容易导致空间流形结构的丧失。

SLLE(Supervised Locally Linear Embedding)是一种较为充分利用分类先验知识的有监督的 LLE 算法^[5], 本文将 SLLE 思想引入 OPRA, 提出了一种新的基于有监督流形学习的正交投影降维方法。基于有监督流形学习的正交投影降维方法用来处理包含多个分离流形的数据, 其中每个流形分别对应各个不同的类别。对于完全分离的流形来说, 一个属于类别 $c'(0 < c' < c)$ 中的样本点 x_i 的近邻点应该从类别 c' 的数据点中寻找。同时, 本文通过人为地增加类间样本之间的距离, 而保持类内样本之间距离来引入更多分类先验信息。样本之间的新距离为

$$D' = D + \alpha \max(D)\gamma, \alpha \in [0, 1] \quad (2)$$

其中, $\max(D)$ 是数据点集的欧氏距离矩阵中的最大值, 如果 x_i 和 x_j 属于同一类, 则 $\gamma_{ij} = 0$, 否则 $\gamma_{ij} = 1$ 。如果 $\alpha = 0$, 则退化为第 2 节介绍的无监督的 OPRA, 简记为 0-OPRA; 如果 $\alpha = 1$, 则为文献[3]提到的有监督的 OPRA(S-OPRA), 如果 α 在 0 和 1 之间变化则可以看成是部分监督的 OPRA, 简记为 α -OPRA, 对于 S-OPRA 来说, 类间数据点中的距离和整个数据集中的最大距离一样大, 这样就意味着样本邻接点只能从相同类别的点集中选择。图 1 给出了 OPRA、S-OPRA 与 α -OPRA 的近邻点计算示例。其中, 近邻点个数 $K=4$ 。

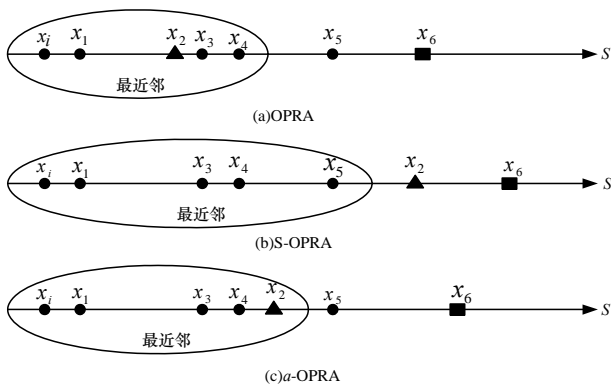


图 1 OPRA, S-OPRA 与 α -OPRA 的近邻点计算示例

在图 1 中, x_i, x_1, x_3, x_4, x_5 属于类 1, 用圆表示; x_2 属于类 2, 用三角形表示; x_6 属于类 3, 用正方形表示。 α -OPRA 加入了额外的参数 α , 通过调节 $\alpha(0 < \alpha < 1)$ 来控制监督的程度, 从而使高维到低维的映射在保留流形某些结构的同时, 也进一步获得了更好的正交投影效果, 它在纯粹的有监督 OPRA 和无监督 OPRA 之间取得了某些折中。

4 结果与分析

图 2(a)为 3D 空间中的 SWISSROLL 流形结构; 根据原始 3D 流形结构随机生成 800 个空间三维点构成测试样本集(图 2(b)), 随机生成函数为

$$X = \{(t \cos(t), ar, t \sin(t))\}, t = \frac{3\pi}{2}(1 + br) \quad (3)$$

其中, $a=21, b=2, r$ 为 0~1 间的随机数; 根据图 2(b)将流形 3D 离散点分为 107 个类别, 并采用不同灰度表示不同类别。表 1 为部分测试样本信息。

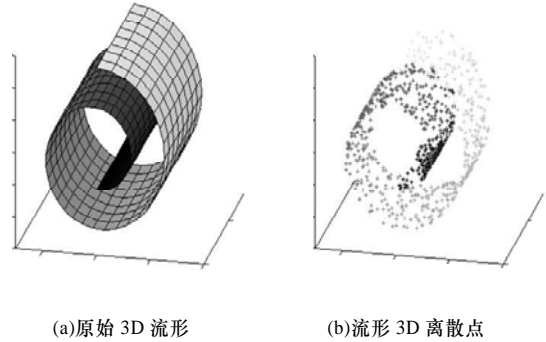


图 2 流形 3D 离散点

表 1 SWISSROLL 流形 3D 离散点部分测试样本信息

编号	X(1)	X(2)	X(3)	类别
1	1.663 4	4.969 6	-11.020 6	1
2	4.326 2	7.744 8	-3.549 7	2
3	2.105 2	7.031 5	-10.985 0	1
4	6.820 5	17.832 7	-9.411 1	3
5	12.497 4	3.968 7	-0.651 0	4
6	3.120 3	11.587 2	-10.835 6	5
7	-7.751 2	11.894 3	-6.511 9	6
8	6.278 0	8.412 8	2.017 1	7
9	-8.119 4	16.277 2	-5.932 5	8
10	12.277 2	2.856 3	3.865 6	9
...

图 3 为不同 α 值 α -OPRA 分别对 3D 空间中 SWISSROLL 流形结构进行 2D 降维后的比较图, 可以看出, α 取 0.1 时(图 3(c)), α -OPRA 效果最好, α 为 0 时(图 3(a)), 是无监督的 OPRA, 效果相对较差, 而 α 大于 0.1 时, 随着 α 的增加, 空间中的部分流形结构逐渐丧失, 效果逐渐变差; α 取 1 时(图 3(f)), 为文献[3]提到的有监督的 OPRA(S-OPRA), 导致空间中的流形结构的丧失, 效果最差。

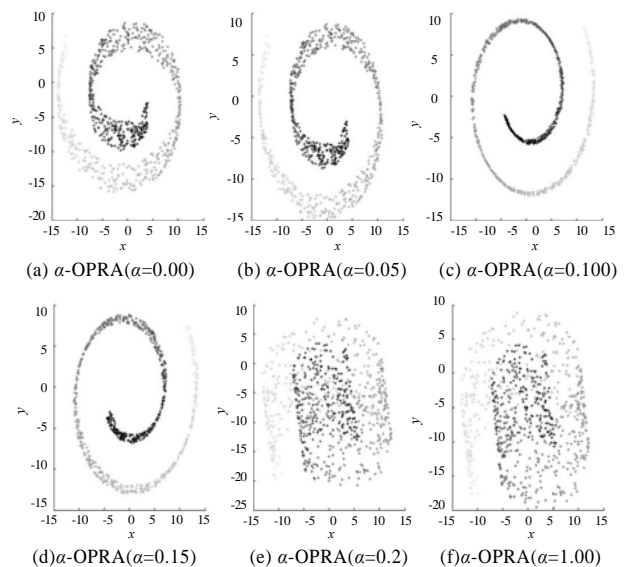


图 3 不同 α 值 α -OPRA 降维效果的比较

(下转第 211 页)