

# 基于核函数的稳健线性嵌入方法

徐雪松<sup>1)</sup> 宋东明<sup>1)</sup> 张 谔<sup>1)</sup> 许满武<sup>2)</sup> 刘凤玉<sup>1)</sup>

<sup>1)</sup>(南京理工大学计算机科学与技术学院, 南京 210094) <sup>2)</sup>(南京大学计算机科学与技术系, 南京 210093)

**摘 要** LLE算法是一种新的非监督学习方法,主要针对非线性降维问题。针对该算法存在的缺点,提出了一种基于核函数的稳健线性嵌入方法,该方法通过引进核函数来优化算法邻域点的求解;在特征空间中,修正权值矩阵 $W$ ,进行降噪处理,经过推导,最终将实际的子空间计算归结为标准的特征值分解问题。采用最小近邻分类器估算识别率。在Yale人脸库以及AT&T人脸库的测试结果表明,在姿态、光照、表情、训练样本数目变化的情况下,改进的算法都具有较好的识别率。

**关键词** 流形学习 高维数据 维数约减 核函数

中图法分类号: TP391.41 文献标识码: A 文章编号: 1006-8961(2009)06-1141-07

## Robust Linear Embedding Based on a Kernel Function

XU Xue-song<sup>1)</sup>, SONG Dong-ming<sup>1)</sup>, ZHANG Xu<sup>1)</sup>, XU Man-wu<sup>2)</sup>, LIU Feng-yu<sup>1)</sup>

<sup>1)</sup>(Department of Computer Science and Technology, Nanjing University of Science and Technology, Nanjing 210094)

<sup>2)</sup>(Department of Computer Science and Technology, Nanjing University, Nanjing 210093)

**Abstract** As a new unsupervised learning method, Local Linear Embedding algorithm (LLE) aims at reducing the nonlinear dimensionality. Since the local linear embedding method has many disadvantages, a new method, namely robust linear embedding method based on a kernel function, is presented to solve this problem. Firstly, the kernel function is utilized to adjust the Euclidean distance between data points, so the new method can improve the performance and the range of application of LLE. Secondly, the new method using the improved  $W$  is selected because it is insensitive to noise. It is shown that the actual computation of the subspace is reduced to a standard eigenvalue problem. The proposed method was tested and evaluated in the Yale face database and AT&T face database. Nearest neighborhood (NN) algorithm was used to construct classifiers. The experimental results showed that the improved algorithm has good performance when pose, lighting condition, face expression and train sample number change.

**Keywords** manifold learning, high dimensional data, dimensionality reduction, kernel function

## 1 引言

人脸图像的维数通常是非常高的,实际上,人脸图像在这样高维空间中的分布很不紧凑,因而不利于分类。并且在计算上的复杂度也非常大。为得到较紧凑分布的人脸图像,Kirby等人 and Turk等人首次把主元分析的子空间思想引入到人脸识别

中<sup>[1-2]</sup>,并获得了较大的成功<sup>[3]</sup>。子空间分析方法随后引起了人们的广泛关注,从而成为了当前人脸识别的主要方法之一。子空间分析的思想是根据一定的性能目标来寻找线性或非线性的空间变换,把原始信号数据压缩到一个低维子空间,使数据在子空间中的分布更加紧凑,为更好地描述数据提供手段,同时计算的复杂度也得到了大大降低<sup>[4-5]</sup>。

由于线性子空间方法不足以有效地描述人脸图

收稿日期:2007-10-08;改回日期:2008-02-22

第一作者简介:徐雪松(1975~),男,南京理工大学计算机应用专业博士研究生。主要研究方向为模式识别、机器学习及信息安全。

E-mail: xxs105@yahoo.cn

像中诸如光照、表情和姿态等复杂的非线性变化,而人脸空间更可能是一个非线性子空间,所以 Roweis 等人于 2000 年提出了 LLE (locally linear embedding) 算法,是一种非线性降维方法<sup>[6]</sup>。LLE 算法可以学习到非线性流形的结构,但是它存在几个问题:(1)该算法要求样本在流形上是稠密采样<sup>[7]</sup>,导致在高维稀疏空间中采用欧氏距离公式失效;(2)在求解模型过程中可能出现的矩阵病态,会导致求解过程对于数据点噪声十分敏感<sup>[8]</sup>,使结果受噪声影响很大;(3)该算法的非线性特性往往使获得非线性映射只能定义于训练集上<sup>[9]</sup>,至于如何将其应用于检验集合尚待进一步研究,即若一个点不属于获得变换的训练集,则获得的变换就不能将这一点相应地映射到低维空间中。因此,本文提出一种改进的 LLE 流形算法即基于核函数的稳健线性嵌入方法。该算法基本思想是在 LLE 流形算法中引入核函数距离测度,以克服算法中由于高维稀疏空间所造成欧氏距离公式的点与点之间距离的对比性差等缺点,并从分析噪声对数据集局部特性的影响入手,通过修正权值矩阵  $\mathbf{W}$ ,可以很好地减弱噪声对降维过程的影响;然后在损失函数中引入一个线性变换矩阵,利用非线性核技巧在维数很高的空间里求解子空间,最终达到最佳降维效果。

## 2 LLE 算法

LLE 算法是一种依赖于局部线性的算法,它认为在局部意义下,数据结构是线性的,或者说局部意义下的点在一个超平面上,主要利用局部的线性来逼近全局的非线性,保持局部的几何结构不变,通过相互重叠的局部邻域来提供整体的信息,从而保持整体的几何性质<sup>[10-11]</sup>。

LLE 算法是映射数据  $A = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$ ,  $\mathbf{a}_i \in \mathbf{R}^D$  到数据集  $B = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n\}$ ,  $\mathbf{b}_i \in \mathbf{R}^d$  ( $D > d$ )。该算法主要包括 3 步:

(1) 对高维空间中的每个样本点  $\mathbf{a}_i$  ( $i = 1, 2, \dots, n$ ), 计算它和其他  $n - 1$  个样本点之间的距离, 根据距离的大小, 选择前  $K$  个与  $\mathbf{a}_i$  ( $i = 1, 2, \dots, n$ ) 最近的点作为其邻近点, 常采用欧氏距离来度量两个点之间的距离, 即  $d_{ij} = \|\mathbf{a}_i - \mathbf{a}_j\|$ ;

(2) 对每个  $\mathbf{a}_i$  ( $i = 1, 2, \dots, n$ ), 找到它的  $K$  个邻近点之后, 计算该点和它的每个邻近点之间的权值  $w_j^{(i)}$ , 即最小化:

$$G(w) = \min \sum_{i=1}^n \left\| \mathbf{a}_i - \sum_{j=1}^K w_j^{(i)} \mathbf{a}_j \right\|^2 \quad (1)$$

式中,  $\sum_{j=1}^K w_j^{(i)} = 1$ , 如果  $\mathbf{a}_j$  ( $j = 1, 2, \dots, K$ ) 不是  $\mathbf{a}_i$  ( $i = 1, 2, \dots, n$ ) 的近邻, 则  $w_j^{(i)} = 0$ ;

(3) 根据高维空间中的样本点  $\mathbf{a}_i$  ( $i = 1, 2, \dots, n$ ) 和它的近邻  $\mathbf{a}_j$  ( $j = 1, 2, \dots, K$ ) 之间的权值  $w_j^{(i)}$  来计算低维嵌入空间中的值  $\mathbf{b}_i$  和  $\mathbf{b}_j$ 。由于在低维空间中尽量保持高维空间中的局部线性结构, 而权值  $w_j^{(i)}$  代表着局部信息, 所以固定权值  $w_j^{(i)}$ , 使下面的损失函数最小化:

$$L(\mathbf{B}) = \min \sum_{i=1}^n \left\| \mathbf{b}_i - \sum_{j=1}^K w_j^{(i)} \mathbf{b}_j \right\|^2 = \text{tr}(\mathbf{B}^T \mathbf{M} \mathbf{B}) \quad (2)$$

式中,  $\mathbf{M} = (\mathbf{I} - \mathbf{W})^T (\mathbf{I} - \mathbf{W})$ ,  $\mathbf{I}$  为单位阵。

要求  $\sum_{i=1}^n \mathbf{b}_i = 0$  且  $\frac{1}{n} \sum_{i=1}^n \mathbf{b}_i \mathbf{b}_i^T = 1$ , 以使  $L(\mathbf{B})$  对

平移、旋转和伸缩变化都具有不变性, 使  $L(\mathbf{B})$  最小化的解为矩阵  $\mathbf{M}$  的最小几个特征值所对应的特征向量构成的矩阵  $\mathbf{B}$ <sup>[12]</sup>。取  $\mathbf{M}$  最小的  $m + 1$  个特征值对应的特征向量, 去掉其中最小的特征值对应的特征向量, 剩余的  $m$  个特征向量组成的矩阵就是低维空间中所得特征向量。

## 3 基于核函数的稳健线性嵌入方法

### 3.1 算法概述

假设输入空间数据集  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ ,  $\mathbf{x}_k \in \mathbf{R}^D$ , 首先, 该输入空间数据集被某种非线性映射  $f$  映射到特征空间 (或是一个 Hilbert 空间) 得到  $f(\mathbf{x}_1)$ ,  $f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)$ 。基于核函数的稳健线性嵌入方法的目的是对特征空间中的数据点 ( $f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)$ ) 进行降维, 把它们映射到  $d$  维空间中形成新的数据点  $Y = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)$ , 其中  $d \ll D$ 。

(1) 在特征空间中进行线性优化操作, 对特征空间中的每个样本点  $f(\mathbf{x}_i)$  ( $i = 1, 2, \dots, n$ ), 计算它和其他  $n - 1$  个样本点之间的距离, 根据核函数距离的大小, 选择前  $K^*$  ( $K^* = 1, 2, \dots, n$ ) 个与  $f(\mathbf{x}_i)$  ( $i = 1, 2, \dots, n$ ) 最近的点作为其邻近点, 距离公式如下节详解, 由参数  $\sigma$  决定  $K^*$  个邻近点紧致性调整;

(2) 对每个  $f(\mathbf{x}_i)$  ( $i = 1, 2, \dots, n$ ), 找到它的  $K^*$  个邻近点之后, 将修正后的权值能量加入极小化目标函数使得

$$Q(\hat{W}) = \min \left\{ \sum_{i=1}^n \left| f(\mathbf{x}_i) - \sum_{j=1}^{k^*} w_j^{(i)} f(\mathbf{x}_j) \right|^2 + \psi \|\mathbf{W}\|^2 \right\} \quad (3)$$

这样最终得到拓扑结构矩阵  $\tilde{W}$ , 式中,  $\sum_{j=1}^{k^*} w_j^{(i)} = 1$ , 如果  $f(\mathbf{x}_j)$  ( $j=1, 2, \dots, k^*$ ) 不是  $f(\mathbf{x}_i)$  ( $i=1, 2, \dots, n$ ) 的近邻, 则  $w_j^{(i)} = 0$ ,  $\Psi$  为系数;

(3) 将线性变换矩阵引入 LLE 算法的损失函数中, 并将其映射至特征空间, 在特征空间中研究损失函数最小化问题;

(4) 分解核矩阵  $k$  并计算矩阵  $M$  和  $Q$ , 且  $K = PAP^T, M = (I - W)(I - W)^T, Q = P^T M P$ , 其中  $I$  表示单位矩阵;

(5) 求解下面的特征值分解问题得到向量  $R_i$ :

$$QR_i = l_i^* R_i \quad (4)$$

$l_i^*$  ( $i=1, 2, \dots, n$ ) 为特征值, 其中  $0 < l_1 < \dots < l_d$  ( $d < n$ );

(6) 计算向量  $\beta_i$  并对其归一化:

$$\beta_i = PA^{-1} R_i$$

$$\beta_i \leftarrow \beta_i / \sqrt{\beta_i^T K \beta_i} \quad (5)$$

(7) 用式(6)计算非线性降维

$$y_i^n = \sum_{j=1}^n \beta_j^n K_{i,j} \quad (6)$$

式中,  $y_i^n$  表示向量  $y_i$  的第  $n$  个分量。

### 3.2 算法的证明及分析

对步骤(1)中所用核距离公式推导如下:

在 Hilbert 空间中, 任意两点  $\mathbf{x}_1, \mathbf{x}_2$  之间的点积为  $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle$ , 定义特征空间中一点的范式为  $\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$ 。由于点积满足对称和双线性性质, 因此, 可以用范式定义特征空间中两点  $\mathbf{x}_1, \mathbf{x}_2$  之间的距离为

$$\|\mathbf{x}_1 - \mathbf{x}_2\| = \sqrt{\langle \mathbf{x}_1 - \mathbf{x}_2, \mathbf{x}_1 - \mathbf{x}_2 \rangle}$$

$$= \sqrt{\langle \mathbf{x}_1, \mathbf{x}_1 \rangle - 2\langle \mathbf{x}_1, \mathbf{x}_2 \rangle + \langle \mathbf{x}_2, \mathbf{x}_2 \rangle} \quad (7)$$

这样可以基于特征映射  $f$  来定义空间数据集中不同数据之间的关系, 定义距离函数  $f_d$  为

$$f_d(\mathbf{d}_1, \mathbf{d}_2) = \|f(\mathbf{d}_1) - f(\mathbf{d}_2)\| = \sqrt{\langle f(\mathbf{d}_1), f(\mathbf{d}_1) \rangle - 2\langle f(\mathbf{d}_1), f(\mathbf{d}_2) \rangle + \langle f(\mathbf{d}_2), f(\mathbf{d}_2) \rangle} \quad (8)$$

但在很多情况下, 由于特征空间维数很高, 很难显示出一个数据映射到特征空间。然而, 本文关心

的仅仅是两个数据在特征空间内的像的点积, 因此, 不必显示计算出的两个像。

定义核函数来计算两个数据特征空间中的像的点积, 它以两个数据点作为输入, 输出是两个数据在特征空间中的像的点积:

$$F_f(\mathbf{d}_1, \mathbf{d}_2) = \langle f(\mathbf{d}_1), f(\mathbf{d}_2) \rangle \quad (9)$$

通过核函数可以重新定义式(7)为

$$f_d(\mathbf{d}_1, \mathbf{d}_2) = \sqrt{\langle F_f(\mathbf{d}_1, \mathbf{d}_1) - 2F_f(\mathbf{d}_1, \mathbf{d}_2) + F_f(\mathbf{d}_2, \mathbf{d}_2) \rangle} \quad (10)$$

可以通过核函数有效地计算特征空间中数据点之间的距离  $f_d$ , 例如高斯核函数:

$$F_{\text{rb}}(\mathbf{d}_1, \mathbf{d}_2) = \exp\left(-\frac{\|\mathbf{d}_1 - \mathbf{d}_2\|^2}{\sigma^2}\right) \quad (11)$$

因此可得:

$$F_{\text{rb}}(\mathbf{d}_1, \mathbf{d}_2) = \exp\left(-\frac{F_f(\mathbf{d}_1, \mathbf{d}_1) - 2F_f(\mathbf{d}_1, \mathbf{d}_2) + F_f(\mathbf{d}_2, \mathbf{d}_2)}{\sigma^2}\right) \quad (12)$$

核函数中的参数  $\sigma$  决定了对优化邻域紧致性的影响, 参数  $\sigma$  的值越小优化邻域区域越紧致, 但并非越小越好, 太小的  $\sigma$  将导致无效的计算结果, 该参数可根据实际情况进行选取。

对步骤(2)的推导如下:

由于  $W$  是由数据集  $f(X)$  解出的, 故当数据集  $X$  受噪声污染时, 计算对噪声非常敏感, 特别是当  $(f(X))^T f(X)$  的特征值较小时, 甚至得不到需要的结果。因此, 权值矩阵  $W$  是否准确将直接影响低维嵌入的最终效果。下面研究当数据集受噪声污染时, 权值  $w_i$  的变化情况。

令  $f(\mathbf{x}_0)$  代表  $f(\mathbf{x}_i)$   $i=1, \dots, n$ ,  $f(\mathbf{x}_i)$  代表相应的真实值,  $U(f(\mathbf{x}_0))$  代表  $f(\mathbf{x}_0)$  的邻域, 设  $f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_k) \in U(f(\mathbf{x}_0))$ , 则

$$f(\mathbf{x}_0) = \sum_{i=1}^{k^*} W^i f(\mathbf{x}_i) \quad \sum_{i=1}^{k^*} W^i = 1$$

令  $\hat{f}(\mathbf{x}_i) = f(\mathbf{x}_i) + d_i$  ( $i=0, 1, 2, \dots, k$ ) 代表相应的噪声影响点, 以及相应的

$$\hat{f}(\mathbf{x}_0) = \sum_{i=1}^{k^*} \hat{W}^i \hat{f}(\mathbf{x}_i) \quad \sum_{i=1}^{k^*} \hat{W}^i = 1$$

若再令

$$f(\mathbf{X}^0) = (f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_k))$$

$$\hat{f}(\mathbf{X}^0) = (\hat{f}(\mathbf{x}_1), \hat{f}(\mathbf{x}_2), \dots, \hat{f}(\mathbf{x}_k)), \text{ 以及}$$

$$W = (W^1, W^2, \dots, W^k)^T, \hat{W} = (\hat{W}^1, \hat{W}^2, \dots,$$

$\widehat{\mathbf{W}}^k)^T$  于是有

$$f(\mathbf{x}_0) = f(\mathbf{X}^0) \mathbf{W}, \widehat{f}(\mathbf{x}_0) = \widehat{f}(\mathbf{X}^0) \widehat{\mathbf{W}}$$

式中,  $\mathbf{X}^0$  代表  $\mathbf{x}_0$  点的邻域矩阵。

**定理** 在上述记号下, 各点噪声之间, 不同真实值之间, 以及  $\widehat{\mathbf{W}}$  与噪声之间是相互独立的, 各点噪声是同均值(0), 同方差的, 并且记  $\delta \mathbf{W} = \widehat{\mathbf{W}} - \mathbf{W}$ , 有如下判别式:

$$E \|\widehat{\mathbf{W}}\|^2 \geq E \|\delta \mathbf{W}\|^2 \frac{\lambda_{\min} l}{k^* (k^* + 1) \rho^2} \quad (13)$$

式中,  $\|\cdot\|$  取为欧几里德范数,  $l = \text{rank}(f(\mathbf{X}^0))$ ,  $\lambda_{\min}$  为  $(f(\mathbf{X}^0))^T f(\mathbf{X}^0)$  的最小非零特征值,  $d^i$  代表  $\mathbf{d}$  (距离向量) 的第  $i$  个分量,

$$\rho^2 = \sum_{i=1}^{k^*} \rho_i^2, \rho_i^2 = \text{Var}(d^i) \quad (i = 1, 2, \dots, k) \quad (14)$$

证明: 由  $\widehat{f}(\mathbf{x}_i) = f(\mathbf{x}_i) + d_i (i = 0, 1, \dots, k)$  可见

$$\widehat{f}(\mathbf{x}_0) = \widehat{f}(\mathbf{X}^0) \widehat{\mathbf{W}} + \sum_{i=1}^{k^*} \widehat{\mathbf{W}}^i d^i \Rightarrow f(\mathbf{x}_0) = f(\mathbf{X}^0) \widehat{\mathbf{W}} +$$

$$\sum_{i=1}^{k^*} \widehat{\mathbf{W}}^i (d^i - d^0)$$

进一步有

$$f(\mathbf{X}^0) (\widehat{\mathbf{W}} - \mathbf{W}) = \sum_{i=1}^{k^*} \widehat{\mathbf{W}}^i (d^0 - d^i),$$

$$\sum_{i=1}^{k^*} \widehat{\mathbf{W}}^i = \sum_{i=1}^{k^*} \mathbf{W}^i = \mathbf{I}$$

$$\text{即 } f(\mathbf{X}^0) \delta \mathbf{W} = \sum_{i=1}^{k^*} \widehat{\mathbf{W}}^i (d^0 - d^i) \quad (15)$$

由于噪声是独立的, 则有

$$\begin{aligned} E(\delta \mathbf{W}^T (f(\mathbf{X}^0))^T f(\mathbf{X}^0) \delta \mathbf{W}) &= E \left[ \sum_{i=1}^{k^*} \widehat{\mathbf{W}}^i (d^0 - d^i)^T \sum_{j=1}^k \widehat{\mathbf{W}}^j (d^0 - d^j) \right] \\ &= \sum_{i=1}^{k^*} E(\mathbf{W}^i)^2 \rho^2 + \sum_{i=1}^{k^*} E(\widehat{\mathbf{W}}^i \widehat{\mathbf{W}}^j) \rho^2 \leq \\ &(k^* + 1) \rho^2 \sum_{i=1}^{k^*} E(\mathbf{W}^i)^2 \\ &= (k^* + 1) \rho^2 E \|\widehat{\mathbf{W}}\|^2 \quad (j = 1, \dots, k) \end{aligned} \quad (16)$$

记  $(f(\mathbf{X}^0))^T f(\mathbf{X}^0)$  的正交分解为

$$(f(\mathbf{X}^0))^T f(\mathbf{X}^0) = \mathbf{S}^T \mathbf{A} \mathbf{S}$$

$$\text{式中, } \mathbf{S} \text{ 为正交阵, } \mathbf{A} = \begin{bmatrix} \lambda_1 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & \lambda_k \end{bmatrix}.$$

由于  $\text{rank}(f(\mathbf{X}^0))^T f(\mathbf{X}^0) = \text{rank}(f(\mathbf{X}^0))$ ,

所以

$$E(\delta \mathbf{W}^T (f(\mathbf{X}^0))^T f(\mathbf{X}^0) \delta \mathbf{W}) = E(\delta \mathbf{W}^T \mathbf{S}^T \mathbf{A} \mathbf{S} \delta \mathbf{W}) \geq$$

$$\lambda_{\min} \sum_{i=1, \lambda_i > 0}^k E(\delta \mathbf{W}^i)^2 \quad (17)$$

通过对  $\mathbf{S}$  进行行初等变换(改变两行的位置), 然后计算相应的式(17), 将所得结果相加可得

$$E(\delta \mathbf{W}^T (f(\mathbf{X}^0))^T f(\mathbf{X}^0) \delta \mathbf{W}) = E(\delta \mathbf{W}^T \mathbf{S}^T \mathbf{A} \mathbf{S} \delta \mathbf{W}) \geq \lambda_{\min} E \|\delta \mathbf{W}\|^2 \frac{1}{k^*} \quad (18)$$

其中,  $\lambda_{\min} = \min_{1 \leq i \leq k^*} \{\lambda_i > 0\}$ 。

结合式(16)可知

$$E \|\delta \mathbf{W}\|^2 \leq \frac{k^* (k^* + 1)}{\lambda_{\min} l} \sum_{i=1}^{k^*} E(\delta \mathbf{W}^i)^2 \quad (19)$$

即

$$E \|\widehat{\mathbf{W}}\|^2 \geq E \|\delta \mathbf{W}\|^2 \frac{\lambda_{\min} l}{k^* (k^* + 1) \rho^2}$$

由上述定理可知, 在邻域大小  $k^*$  已知情况下,  $\mathbf{W}$  的误差主要由 3 个因素决定: (1) 噪声的影响, 即  $\mathbf{d}$  的大小; (2) 邻域的影响, 即  $\lambda_{\min}$  和秩  $l$  的大小; (3) 权值能量的  $\|\widehat{\mathbf{W}}\|$  的大小。因素(1)的影响来自于数据采样本身, 因素(2)的影响可通过步骤 1 改变邻域来调整, 本文主要通过修正因素(3)即是通过减少权值能量达到减弱噪声影响的目的。

步骤(3)至步骤(7)推导如下:

LLE 算法中式(2)

$$L(\mathbf{B}) = \min \sum_{i=1}^n \left\| \mathbf{b}_i - \sum_{j=1}^K w_j^{(i)} \mathbf{b}_j \right\|^2 = \text{tr}(\mathbf{B}^T \mathbf{M} \mathbf{B})$$

最小化损失函数的约束项为  $\mathbf{B} \mathbf{B}^T = \mathbf{I}$ 。其中  $\text{tr}(\cdot)$  代表对矩阵求迹的算子。然而该最小化损失函数仅仅依赖于权值  $w$ , 其输出  $\mathbf{b}$  没有显示地与输入点  $\mathbf{a}$  相关联, 也就不能产生便于应用的、可延展的变换函数。为了克服这一问题, 把下面的线性变换矩阵引入最小化损失函数中:

$$\mathbf{B} = \widehat{\mathbf{S}}^T \mathbf{X} \quad \text{或} \quad \mathbf{b}_i = \widehat{\mathbf{S}}^T \mathbf{X}_i \quad (20)$$

其中,  $\mathbf{S} = \{\mathbf{s}_1, \dots, \mathbf{s}_d\}$  于是有

$$\begin{aligned} L(\mathbf{B}) &= \text{tr}(\mathbf{B}^T \mathbf{M} \mathbf{B}) = \text{tr}[(\widehat{\mathbf{S}}^T \mathbf{X}) \mathbf{M} (\widehat{\mathbf{S}}^T \mathbf{X})^T] \\ &= \text{tr}[\widehat{\mathbf{S}}^T (\mathbf{X} \mathbf{M} \mathbf{X}^T) \widehat{\mathbf{S}}] \end{aligned} \quad (21)$$

st

$$\widehat{\mathbf{S}}^T \mathbf{M} (\widehat{\mathbf{S}}^T \mathbf{X})^T = \widehat{\mathbf{S}}^T (\mathbf{X} \mathbf{X}^T) \widehat{\mathbf{S}} = \mathbf{I} \quad (22)$$

利用 Lagrange 乘子求解上述受约束的最小化问题:

$$L(\widehat{\mathbf{S}}) = \widehat{\mathbf{S}}^T (\mathbf{X} \mathbf{M} \mathbf{X}^T) \widehat{\mathbf{S}} + l_i^* (\mathbf{I} - \widehat{\mathbf{S}}^T \mathbf{X} \mathbf{X}^T \widehat{\mathbf{S}}) \quad (23)$$

使得对于  $\widehat{\mathbf{S}}$  的梯度为 0, 得到

$$(XMX^T)\hat{S} = l_i^*(XX^T)\hat{S} \quad (24)$$

或者

$$(XMX^T)s = l_i^*(XX^T)s \quad (25)$$

现在开始在高维特征空间  $F$  中研究上述最小化问题,在式(25)中引入非线性映射函数  $f$ :

$$[f(X)M(f(X))^T]s = l_i^*[f(X)(f(X))^T]s \quad (26)$$

因为特征向量可以用  $F$  中元素的线性组合表示,所以存在系数  $\beta_i (i=1, \dots, n)$  使得:

$$s = \sum_{i=1}^n \beta_i f(x_i) \quad (27)$$

把式(27)代入式(26)中并且两边乘以  $[f(x_j)]^T$  得到:

$$\begin{aligned} & [f(x_j) \cdot f(x_1), \dots, f(x_j) \cdot f(x_n)] M \sum_{i=1}^n \beta_i \cdot \\ & \begin{bmatrix} f(x_1) \cdot f(x_1) \\ \vdots \\ f(x_n) \cdot f(x_1) \end{bmatrix} = l_i^* [f(x_j) \cdot f(x_1), \dots, f(x_j) \cdot f(x_n)] \cdot \\ & \sum_{i=1}^n \beta_i \cdot \begin{bmatrix} f(x_1) \cdot f(x_1) \\ \vdots \\ f(x_n) \cdot f(x_1) \end{bmatrix} \end{aligned} \quad (28)$$

定义核矩阵  $K_{n \times n}$  (它的元素为)  $K_{i,j} = (f(x_i))^T f(x_j) = f(x_i) f(x_j)$  和向量  $\beta = [\beta_1, \dots, \beta_n]$  上式可写成更加紧凑的形式:

$$K M K \beta = l_i^* \tilde{K} \beta \quad (29)$$

上式计算可归结为广义特征值分解问题,如果矩阵  $\tilde{K}$  可逆,则式(29)的广义特征值分解可以转换成标准特征值分解问题。只需在公式两边同时乘以逆矩阵  $(\tilde{K})^{-1}$ 。但实际情况是  $\tilde{K}$  常常为奇异矩阵,为了避免此问题对式(29)继续推导:

在式(29)两边同时乘  $\beta^T$ , 得到

$$\beta^T K M K \beta = l_i^* \beta^T \tilde{K} \beta \quad (30)$$

对  $K$  实施特征值分解:

$$k = P \Lambda P^T \quad (31)$$

其中,  $P$  是正交矩阵而  $\Lambda$  是对角矩阵,把式(31)代入式(30),有

$$(\Lambda P^T \beta)^T P^T M P (\Lambda P^T \beta) = l_i^* (\Lambda P^T \beta)^T (\Lambda P^T \beta) \quad (32)$$

其中,应用了  $P P^T = I$  这一关系。

进一步,令  $R = \Lambda P^T \beta$  从而  $\beta = P \Lambda^{-1} R$ , 式(32)

可简化为  $R^T P^T M P R = l_i^* R^T R$  因此,下式成立

$$P^T M P R = l_i^* R \quad (33)$$

最后定义  $Q = P^T M P R$ , 于是有

$$Q R = l_i^* R \quad (34)$$

注意到式(34)就是步骤(5)中的式(4)。同时在步骤(6)中对  $\beta_i$  的归一化,用来保证与之对应的向量  $R_i$  在  $F$  空间里归一的。联立式(20)和(27)就可以得到步骤(7)中的式(6)。至此,完成了所有推导。

### 4 算法的实现及其实验结果的评估

将本文提出的基于核函数的稳健线性嵌入方法,在 Yale 人脸库以及 AT&T 人脸库上分别进行了测试。Yale 人脸库用来测试在表情和光照变化的情况下算法的性能。AT&T 人脸库用来测试训练样本数目变化时算法的识别情况。图 1 是两个人脸库的图片示例。



图1 AT&T及Yale人脸库的人脸照片  
Fig.1 Sample images for one subject of AT&T and Yale database

Yale 人脸库包括 15 个人,每人 11 张  $195 \times 231$  的灰度照片。这些照片面部表情、面部细节、光照等都有着不同程度的变化。在预处理阶段,将 Yale 人脸库中所有图片进行裁剪,归一化成  $100 \times 100$  的标准人脸图像。在 Yale 库中,将每人的前 5 幅人脸图像样本作为训练样本,后 6 幅图像作为测试样本,这样训练样本总数为 75 个,测试样本总数为 90 个。AT&T 人脸库包括 40 个人,每人 10 张  $112 \times 92$  的灰度照片。照片中人脸的位置和尺度都有不同程度的变化。在 AT&T 人脸库中每人随机选取 5 张照片用于训练,其余的用于测试,一共有 200 个训练样本和 200 个测试样本。我们将本文提出的算法分别与 Eigenface (PCA) 和 LLE 算法进行比较,采用最小近邻分类器估算识别率。

在两个人脸数据库上的对比实验中采用“留一法”验证策略,该验证策略实际上是交叉验证的极端情况。这种方法的基本思路是:假定有  $n$  个初始样本,每次拿出其中一个样本做测试,其余  $n - 1$  个做训练样本,这样的过程重复  $n$  次,结果每个样本都

有一次作为测试样本。其中本文用  $n - 1$  个样本构建一个分类器,再用另外一个样本做测试,得到一个准确度,0 或者 100%,因为这样的过程做了  $n$  次,因此最后的准确度值为  $n$  次的平均值。这个方法有两个明显的好处:第一,它尽可能的使用了所有能用的数据来进行训练;第二,这个方法具有确定性,不需要像交叉验证那样要用随机函数来选择训练集与测试集的划分。

对于 RLEBKF 算法采用的高斯核

$$F_{cb}(d_1, d_2) = \exp(-\|d_1 - d_2\|^2 / \sigma^2)$$

核函数中的参数  $\sigma$  决定了对优化邻域紧致性的影响。表 1 表示在 AT&T 人脸库及 Yale 人脸库上, RLEBKF 算法选取不同参数  $\sigma^2$  的值的识别精度。

表 1 核函数中的  $\sigma$  与优化区域正确识别率之间的关系

Tab.1 Relationship between  $\sigma$  of kernel function and right recognition rates in the optimize region

参数 $\sigma^2$	0.5	0.3	0.1	0.05	0.03
AT&T 人脸库 正确识别率	65	72	79	85	87.5
Yale 人脸库 正确识别率	68	74	80	84.33	87

从表中可以看出,随着  $\sigma^2$  的减小,优化区域识别精度逐渐增大,当  $\sigma^2 = 0.03$  时,识别精度已经达到 87% 以上。

在实验中可以看出,原 LLE 算法在保持原数据集特性的同时,由于数据集受到了噪声的污染,使得结果产生了一定的厚度,即将本征维数得到了增加,这显然有悖于处理的初衷。相反,本文算法则在很好地保持数据集特性的基础上,抑制了噪声的影响,可更加清晰地揭示数据集的本征结构。

表 2 表示在 Yale 人脸库及 AT&T 人脸库中,比较了 3 种不同识别方法可以达到的最佳识别率  $R_{max}$  和所需的特征维数  $D_{im}$ 。可以看出,在 Yale 人脸库及 AT&T 人脸库中 PCA 的识别率最低;LLE 虽然识别率较高,但是特征维数分别达  $42 \times 10$  和  $54 \times 14$ ;而本文方法分别只需要  $9 \times 9$  和  $12 \times 12$  的特征维数,都仅为 LLE 算法的约 1/5,而且识别率最高。

表 3 的(a)和(b)比较了在保证最佳识别率的前提下,不同识别方法所需要的计算时间(整个实验运行在主频为 P4 2.4 GHz、内存为 512 MB、80 GB 硬盘、操作系统为 Windows XP sp2 的主机上面),其中,  $t_E$  为特征提取时间,  $t_C$  为分类时间,  $t_T$  为总时间。

在 Yale 人脸库中可以看出,PCA 在特征提取上耗费的时间最多,分别为 LLE 的 50 倍和本文方法的 37 倍。虽然 PCA 特征维数大,但是其分类时间小于后 2 种方法。同时,虽然 LLE 算法提取特征耗费的时间比 RLEBKF 短,但是由于其特征维数大,计算复杂度高,分类耗费的时间为本文提出的方法的 1.8 倍,所以总时间还是比本文提出的方法多。同样,在 AT&T 人脸库中也有类似结果。

表 2 3 种识别方法的最佳识别率和相应特征维数的比较

Tab.2 Comparison of maximal recognition rate and corresponding feature dimension of three recognition methods

识别方法	$R_{max}(\%)$		$D_{im}$		“留一法” $R_{max}(\%)$	
	Yale	AT&T	Yale	AT&T	Yale (160/165)	AT&T (398/400)
PCA	65.0	74.2	100	150	86.2	89.3
LLE 算法	75.0	86.5	$42 \times 10$	$54 \times 14$	92.2	96.5
本文方法	93.5	95.6	$9 \times 9$	$12 \times 12$	98.75	99.5

表 3 在最佳识别率下,3 种识别方法在不同人脸库中的计算时间

Tab.3 Comparison of CPU time of three recognition methods under maximal recognition rate on face database

识别方法	识别方法	单位:s		
		$t_E$	$t_C$	$t_T$
Yale	PCA	46.35	0.19	46.65
	LLE 算法	0.91	4.78	5.69
	RLEBKF	1.24	2.71	3.95
AT&T	PCA	54.25	0.68	54.93
	LLE 算法	2.28	6.54	8.82
	本文方法	3.68	5.35	9.03

## 5 结 论

LLE 算法非线性局部嵌入方法在生物数据分类、文字识别、脸谱图像处理中获得了一定的效果。但其存在固有的缺点,即该算法要求样本在流形上是稠密采样,对样本中的噪声很敏感,并且不能产生便于应用的、可延展的变换函数。针对其缺点,在 LLE 流形算法中引入核函数距离测度,以克服算法中由于高维稀疏空间所造成欧氏距离公式的点与点之间距离的对比性差等缺点,并从分析噪声对数据

集局部特性的影响入手,通过修正权值矩阵  $W$ ,可以很好地减弱噪声对降维过程的影响;然后在损失函数中引入一个线性变换矩阵,利用非线性核技巧在维数很高的空间里求解子空间,最终达到最佳降维效果。理论分析及在 Yale 人脸库以及 AT&T 人脸库的测试结果表明,基于核函数的稳健线性嵌入方法性能优于 LLE 算法,和其他子空间方法相比,也具有更好的识别能力。

### 参考文献 (References)

- 1 Kirby M, Sirovich L. Application of the karhunen2loeve procedure for the characterization of human faces[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 1990, **12**(1): 103-108.
- 2 Turk M, Pentland A. Eigenfaces for recognition[J]. Journal of Cognitive Neuroscience, 1991, **3**(1): 72-86.
- 3 Phillips P J, Moon H, Rizvi S, et al. The feret evaluation methodology for face recognition algorithms[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2000, **22**(10): 1090-1104.
- 4 Turk M, Pentland A. Eigenfaces for recognition[J]. Journal of Cognitive Neuroscience, 1991, **3**(1): 71-86.
- 5 Zhao W, Chellappa R, Rosenfeld A, et al. Face recognition: A literature survey [R]. CS-Tech Report 4167, University of Maryland, 2000.
- 6 Roweis S T, Saul L K. Nonlinear dimensionality reduction by locally linear embedding [J]. Science, 2000, **290**(22): 2323-2326.
- 7 Xu Rong, Jiang Feng, Yao Hong-xun. Overview of manifold learning [J]. CAAI Transactions on Intelligent Systems, 2006, **1**(1): 44-51. [徐蓉,姜峰,姚鸿勋. 流形学习概述[J]. 智能系统学报, 2006, **1**(1): 44-51.]
- 8 Zhao Lian-wei, Luo Si-wei, Zhao Yan-chang, et al. Study on the low-dimensional embedding and the embedding dimensionality of manifold of high-dimensional data[J]. Journal of Software, 2005, **16**(8): 1423-1430. [赵连伟,罗四维,赵艳敬等. 高维数据的低维嵌入及嵌入维数研究[J]. 软件学报, 2005, **16**(8): 1423-1430.]
- 9 Bengio Y, Palment J, Vincent P, et al. Out-of-sample extensions for LLE, isomap, MDS, eigenmaps, and special clustering[A]. Neural Information Processing Systems [C], Cambridge, MA: MIT Press, 2003: 1238-1248.
- 10 Tenenbaum J B, De Silva V, Langford J C. A global geometric framework for nonlinear dimensionality reduction [J]. Science, 2000, **290**(5500): 2319-2323.
- 11 Shashua A, Levin A, Avidan S. Manifold pursuit: A new approach to appearance based recognition [A]. In: Proceeding of 16th International Conference on Pattern Recognition [C], Quebec City, Canada, 2002: 590-594.
- 12 Weng Shi-feng, Zhang Chang-shui, Zhang Xue-gong. Nonlinear dimensionality reduction in the analysis of high dimensional medical data[J]. Journal of Tsinghua University Sci&Tech, 2004, **44**(4): 485-488.