基于特征子空间邻域的局部保持流形学习算法*

王 娜,李 霞,刘国胜

(深圳大学 信息工程学院 电子系, 广东 深圳 518060)

摘 要:局部保持流形学习算法通过保持局部邻域特性来挖掘隐藏在高维数据中的内在流形结构。然而,对于缺乏足够训练样本的高维数据集,或者高维数据集存在非线性结构和高维数据特征中存在冗余、干扰特征,使得在原特征空间中利用欧式距离定义的邻域关系并不能真实反映数据的内在流形结构,从而影响算法的性能。提出利用正约束寻找特征子空间的方法,使得在此子空间中更多的同类样本紧聚,并进一步在该子空间中构建邻域关系来挖掘高维数据的内在流形,形成基于特征子空间邻域特性的局部保持流形学习算法(NFS-LPP 和 NFS-NPE)。它们在一定程度上克服了高维小样本数据集难以正确挖掘内在流形结构的问题,在 Yale 和 ORL 人脸库上的分类和聚类实验验证了其有效性。

关键词: 正约束; 特征子空间; 局部保持; 流形学习

中图分类号: TP391 文献标志码: A 文章编号: 1001-3695(2012)04-1318-04

doi:10.3969/j.issn.1001-3695.2012.04.032

Locality preserving manifold learning algorithm based on neighborhood in feature space

WANG Na, LI Xia, LIU Guo-sheng

(Dept. of Electronics, College of Information Engineering, Shenzhen University, Shenzhen Guangdong 518060, China)

Abstract: Locality preserving manifold learning algorithms always discover intrinsic manifold in high-dimensional data by preserving locality neighborhood structures. However, for high-dimensional data with non-enough training samples, or with nonlinear structure and redundant or interrupted features, it is difficult to directly estimate real neighbor relation defined by Euclidean distance in original feature space. This paper proposed a novel method to find a feature subspace best suited to representing neighborhood relation using positive constraints. In this subspace more inner-class samples come together. Further, constructed neighborhood graph in this subspace to discover intrinsic manifold in high-dimensional data, which caused novel locality preserving manifold learning algorithms called NFS-LPP and NFS-NPE. Experimental results on Yale and ORL face database verify their effectiveness.

Key words: positive constraints; feature subspace; locality preserving; manifold learning

0 引言

随着信息时代的到来,数据的高维性、计算复杂性、冗余性以及其内在规律难以发掘的特点,使得高维数据降维成为一种必要的数据处理技术,现已广泛应用于图像分析、计算机视觉、Web 信息检索、金融数据分析等各大领域。降维方法可以分为线性和非线性降维。线性降维方法通常假设数据存在于全局线性的结构当中,即数据集各特征之间是相互独立的,典型的线性降维方法有主成分分析(principle components analysis, PCA)^[1]和线性判别分析(linear discriminate analysis, LDA)^[2]等。当数据确实存在于全局线性中时,线性降维是非常有效的。然而现实生活中的高维数据往往具有非线性结构,比如人脸图像,其特征是由姿态、表情、光照强度等因素共同决定的,并不是这些因素的简单线性组合而成,采用线性降维方法难以挖掘出高维数据的内在特性。如何从非线性数据中发掘其内在规律显得尤为重要,近年来流形学习成为解决这一问题的研

究热点。流形学习的主要目标是发掘嵌入高维数据空间中的 低维光滑流形,流形学习算法可粗略地分为基于全局保持和局 部保持两大类。全局保持算法主要是保持高维数据的全局特 性,构建度量每一点与其他数据点关系的全局度量矩阵,并将 其转换为内积矩阵,利用内积矩阵的特征分解获取高维数据的 低维表示,典型算法包括等度规映射(isometric mapping, ISO-MAP)[3]、最大方差展开(maximum variance unfolding, MVU)[4] 等。局部保持算法则是假设高维空间数据在局部是线性的,且 具有欧式空间特性,通过构建局部近邻图来保持局部邻域特 性,主要算法包括局部线性嵌入(locally linear embedding, LLE)^[5]、拉普拉斯特征映射(Laplacian eigenmaps, LE)^[6]等。 这些算法在高维和低维空间之间建立定义在训练样本上的隐 式非线性映射,不适用于测试数据。线性化的流形学习算法很 好地解决了这一问题,如局部保持投影(locality preserving projection, LPP)[7]和局部近邻保持投影(neighborhood preserving embedding, NPE)[8] 算法等,它们分别是 LE 和 LLE 算法的线性

收稿日期: 2011-08-31; **修回日期:** 2011-10-17 **基金项目:** 国家自然科学基金资助项目(60902069,61171124);广东省自然科学基金资助项目(9151806001000025)

作者简介:王娜(1977-),女,河北保定人,教授,博士,主要研究方向为机器学习、模式识别(wangna@szu.edu.cn);李霞(1968-),女,四川乐山人,教授,博导,主要研究方向为智能优化、模式识别;刘国胜(1986-),男,湖南长沙人,硕士,主要研究方向为流形学习.

化推广,但它们不对高维数据集作全局线性假设,只期望保持邻域关系对数据集进行降维,在高维训练样本和其低维嵌入之间求得线性变换矩阵,通过此矩阵得到高维新增数据的低维嵌入。局部近邻图的构建直接影响着局部保持流形学习算法性能,它们大多通过在原特征空间中利用欧式距离来衡量样本之间的相似性。然而由于高维数据可能存在非线性结构,以及高维空间的稀疏性,在原空间中通过欧式距离找寻的近邻未必来自于同类,导致因异类近邻关系的保持而影响算法分类性能。

监督信息主要包括类别标签和成对约束两种形式,是另一种监督信息形式。成对约束可以分为正约束和负约束,具有正约束关系的两个样本点属于同一类;反之,具有负约束关系的样本不属于同一类。很显然,成对约束是一种更为一般、更容易获取的监督信息,可以直接从类别标签中获取。本文提出利用正约束寻找特征子空间,使得在此子空间中更多的同类样本紧聚,进而利用子空间中的邻域关系挖掘原空间中的局部邻域结构,提出了基于特征子空间邻域特性的局部保持流形学习算法(NFS-LPP和NFS-NPE)。

1 局部保持流形学习算法

流形学习假定高维数据集位于或近似位于一个嵌入在高 维欧式空间的内在低维流形上。流形每点处都存在一个局部 邻域结构同胚于欧式空间的子集,故局部近邻数据间就具有近 似欧式几何性质。局部保持的流形学习算法通过在高维空间 中建立局部模型来刻画局部几何特性,然后将这种局部几何特 性保持到低维空间,通过局部几何共性的保持和局部邻域的交 叠覆盖来反映整个嵌入在高维空间的内在低维流形。

局部保持的流形学习算法大致通过局部近邻图的构建、局部几何特性的刻画和保持局部几何特性求解全局低维坐标三个步骤来完成。在局部近邻图构建时,所有局部保持流形学习算法都是一致的,都通过寻找每个样本的近邻来构建邻域成员矩阵 Nb,它是一个 $n \times n$ 非对称的方阵,若寻找 x_i 的 k 个近邻,那么邻域成员矩阵 Nb 的第 i 行将有 k 个为 1 的元素。

$$Nb_{ij} = \begin{cases} 1 & x_i \text{ 和 } x_j \text{ 具有邻域关系} \\ 0 & \text{其他} \end{cases}$$
 (1)

不同局部几何特性的刻画产生了不同的局部保持流形学习算法,LPP 和 NPE 是两种典型的局部保持流形学习算法。

在 LPP 算法中,局部几何特性通过局部权值矩阵 W来刻画。W是一个 $n \times n$ 方阵,它描述所有样本两两之间的关系,它对邻域成员矩阵 Nb 中的非零成员给予重新定义。定义权值的方法有多种,常用的有二值法、热核法等。二值法给予所有近邻样本间的权值都赋为 1,将近邻关系在同一尺度下进行保持,这时 W与 Nb 相等;热核法则根据近邻样本间的距离和核宽 σ 来共同定义近邻间的权值,使得近邻关系按距离层次保持。通过最小化目标函数式(2)可以使高维空间的近邻样本根据近邻权值的定义来保持局部邻域的几何特性。

$$\min \sum_{ij}^{n} \|\alpha^{\mathsf{T}} x_i - \alpha^{\mathsf{T}} x_j\|^2 W_{ij}$$
 (2)

其中:n 为训练样本数量, α 为高维样本 x_i 和其低维表示 y_i 间的变换矩阵。

在 NPE 算法中,局部几何特性则通过局部重建权值矩阵来刻画。对于高维空间中的每一个样本 x_i 都通过它的局部邻域成员(Nb中第 i 行的非零成员)来重新表示,局部重建表示

系数通过在 $\sum_{i=1}^{k} W_{ij} = 1$ 的约束下最小化目标函数式(3)来求取。

$$\min \sum_{i=1}^{n} \|x_{i} - \sum_{i=1}^{n} x_{i} W_{ij}\|^{2}$$
 (3)

在求得每一个样本的重建表示系数之后,保持重建系数矩阵不变,使得在低维空间中样本与其邻域间具有相同的表示关系,由此可以使高维空间的局部邻域结构在低维空间中仍然得以保持。

然而,对于局部保持的流形学习算法,要求高维空间必须 是密集采样的,即需要大量的训练样本才估计出内在流形。现 实中,某些高维数据要获得足够多的训练样本是很困难的,比 如在人脸识别中,针对每一个人,要获得足够多的在不同环境 下的人脸图像来训练具有强烈区分特性的人脸模型就非常困 难。另外,人脸数据的高维性,使得样本间非常稀疏,再加上光 照、表情等因素的影响,使得人脸图像存在非线性结构,这些就 导致了在原特征空间直接利用欧氏距离构建的邻域图中,同类 样本较少。而从分类角度上来说,在局部保持的流形学习算法 中,期望保持同类样本的邻域结构而尽可能减少异类样本邻域 关系的保持来提高算法性能。为此,本文提出了一种利用正约 束构造特征子空间的方法(即在原始特征集合中,通过正约束 选择出一组特征子集,使得在该特征子空间中具有正约束关系 的样本相距较近的同时,也使得更多的同类样本紧聚),从而 克服因数据的高维性和非线性,或训练样本难以获取等问题而 导致的在原空间中利用欧式距离构造的局部邻域结构难以正 确挖掘数据内在类别结构的问题。

2 基于特征子空间邻域的局部保持流形学习算法

2.1 特征子空间的构造

设训练样本 $X = (x_1, x_2, \cdots, x_n) \in \mathbb{R}^D$,定义 x_i 为样本点 x_i 的第 t 维特征, $\mu_t = \frac{1}{n} \sum_{i=1}^{n} x_i$ 表示样本在第 t 维上的均值,ML 为正约束对集。定义如下计算方式来对每一维特征作出评价,获得一个评价数组 $\{V_t\}$, $t = 1, 2, \cdots, D$ 。

$$V_{t} = \frac{\sum_{(x_{i}, x_{j}) \in ML} \|x_{it} - x_{jt}\|^{2}}{\frac{1}{n} \sum_{i}^{n} \|x_{it} - \mu_{t}\|^{2}}$$
(4)

从式(4)可以看出,V,表示所有正约束关系的点在第t维特征上的距离之和与对应维度上方差的比值。显然,V,越小,表明在第t维特征上具有正约束关系的点相距越近。将具有较大V,值的维度称之为违反正约束关系的特征维。

正约束对属于同一类,理应相距很近,之所以有些正约束样本相距较远,主要是因为高维数据中存在一些冗余和干扰特征或度量方式的局限性。通过式(4)可以计算出所有维特征的评价数组 $\{V_i\}$,并对其进行升序排列,去除一定数量具有较大 V_i 值的违反特征能排除干扰和冗余特征的影响,降低度量方式的局限性,使得更多的同类样本相距较近。

以 Yale 和 ORL 人脸库为例, Yale 人脸库包含 15 个人在不同光照、表情、姿态下的 165 张人脸图像,每人 11 张,将每张人脸图像规格成 24×24 大小,那么这 165 张人脸图像就相当于576 维空间中的 165 个样本点; ORL 人脸库包含 40 个人在不同表情、是否佩戴眼镜下的 400 张人脸图像,每人 10 张,将人脸图像规格成 32×32 大小,那么它们就相当于 1 024 维空间中的 400 个样本点。在 Yale 人脸库中,对于每一个样本与其

同类的样本数为10,利用欧式距离寻找每一个样本点的10邻 域,则所有样本的邻域点之和为1650,其中同类样本数仅为 1 146(相当于图 1(a)中选择特征比例为 1 时);同样在 ORL 人 脸库中,对于每一个样本与其同类的样本数为9,利用欧式距 离寻找每一个样本点的9邻域,那么所有样本的邻域点之和为 3600,其中同类样本的总数仅为2642(相当于图1(b)中选择 特征比例为1时)。随机产生不同数量(20,30,40)的正约束 对,在Yale和ORL人脸库中利用这些正约束选择出不同比例 的特征来构造特征子空间,统计在该特征子空间下近邻图中同 类样本的数量。从图1中可以看出,无论是在多少对正约束和 保留多少比例特征的情况下,所有样本的近邻中同类样本的数 量都多于在原空间下的情况。而且还可以看出,在 Yale 人脸 库中选择保留 20%~30% 左右的特征和在 ORL 人脸库中选择 保留30%~50%的特征来构建特征子空间时,子空间中的邻 域中来自同类样本的数量较多,远大于在原空间中的情况。可 见,采用这种方法不仅能使得具有正约束关系的点在特征子空 间相距较近,而且使部分原空间相距较远的无正约束标记的同 类样本相距较近。

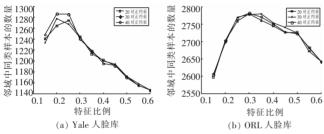


图 1 在不同数量正约束下选择出不同比例特征构建的 特征子空间的邻域中同类样本数量比较

2.2 基于特征子空间邻域的局部保持流形学习算法

在构建 LPP 和 NPE 算法的局部近邻图时,通过正约束选 择特征来构造特征子空间,然后在该子空间中采用欧式距离来 构建邻域图,得到新的邻域成员矩阵 Nb,并将其应用到 LPP 和NPE算法中。虽然在特征子空间中能使得具有正约束关系 的样本在相距较近的同时也能拉近部分在原空间中相距较远 而未被标注成正约束关系的同类样本,但不难预料,当正约束 达到一定数量后,根据式(4)选择保留的特征可能趋于一致 (不会因为正约束数量的增加而使选择出来的特征存在差 异),由此构造的特征子空间则会趋于稳定,以至于算法的性 能得不到进一步的改善。另外,在特征子空间构造的邻域中, 并不能确保包含所有的正约束样本,使得某些正约束样本映射 到低维空间仍然相距较远而影响分类性能,从而使得正约束在 提高算法的分类性能上没有得到充分利用。为充分利用已知 先验正约束信息来提高算法的判别性能,进一步利用正约束来 修正 LPP 算法的权值矩阵和 NPE 算法中的邻域成员,形成了 基于特征子空间邻域的流形学习算法 NFS-LPP 和 NFS-NPE。

在 NFS-LPP 算法中,局部权值矩阵的修正方法如式(5)所示:

$$\widetilde{W} = \overline{W}_{ii} + ML \tag{5}$$

其中: W 是基于特征子空间邻域成员矩阵 Nb 的权值矩阵; ML 为正约束矩阵, 它是一个对称的方阵, ML_{ij} 表示样本 x_i 和 x_i 是 否有正约束, 若有,则其值为 1,否则为 0。

通过修正权值矩阵,可以增大具有正约束关系对的权值, 从而保证所有具有正约束关系的样本在低维空间相距较近。

在 NFS-NPE 中重构各个样本时,尽可能利用与之同类的

样本来重构,对基于特征子空间的邻域成员矩阵 Nb 进行修正,以减小重构误差。将与之有正约束关系的样本纳入到其子空间邻域中,而将相距相对较远的样本从其邻域中剔除。比如,若在子空间中寻找 x_i 的 k 个近邻点,而且已知两个邻域外的样本与 x_i 具有正约束关系,那么就将这两个样本纳入到 x_i 的邻域中,而将与 x_i 相距较远的两个样本从其邻域中移除。

通过这种修正邻域成员的方法获得修正的邻域成员矩阵 \widetilde{Nb} ,

并采用式(3)求得修正的局部重建系数矩阵 \tilde{W} 来尽可能使得所有具有正约束关系的样本映射到低维空间相距很近。算法的具体描述如表 1 所示。

表1 算法描述

NFS-LPP 算法 NFS-NPE 算法 1. 输入训练样本集 $X = (x_1, x_2, \cdots, x_n) \in \mathbb{R}^D$,根据已知先验正约束信息构造正约束矩阵 ML 2. 利用式(4)对 D 维特征作出评价,根据分数进行升序排列,选择分数较低的 D_1 个特征作为特征选择的结果

- 3. 根据选择的 D_1 个特征构造特征子空间 $\bar{X}=(\bar{x}_1,\bar{x}_2,\cdots,\bar{x}_n)\in R^{D_1}$,并在特征子空间中寻找每个样本的 k-近邻,构建邻域成员矩阵 $\overline{N}b$
- 4. 根据 $\overline{N}b$ 定义权值矩阵 \overline{W} , 根据式 (5) 求得修正权值矩阵 \widetilde{W}
- 5. 构建目标函数 $\min = \sum_{ij}^{n} \| \alpha^{T} x_i \alpha^{T} \|$ $\| x_i \|^2$,约束条件为: $\alpha^{T} X X^{T} \alpha = I$
- 6. 将帯约束的目标函数转换为如下广义特征方程: $XLX^{T}\alpha = \lambda XX^{T}\alpha$, 其
 - $\ddagger L = D \widetilde{W}, D_{ii} = \sum_{i,j} \widetilde{W}_{ij}$
- 4. 利用正约束对 $\overline{N}b$ 进行修正,得到修正的邻域成员矩阵 $\widetilde{N}b$,通过式(3)求得修正的局部重建系数矩阵 \widetilde{W} ;
- 5. 构建目标函数 $\min = \sum_{i}^{n} \| \alpha^{T} x_{i} \sum_{j}^{n} \alpha^{T} x_{j} \tilde{W}_{ij} \|^{2}$, 约 束 条 件 为:
- $\alpha^{T}XX^{T}\alpha = I$ 6. 将带约束的目标函数转换为如下广 义特征方程: $XMX^{T}\alpha = \lambda XX^{T}\alpha$, 其 中 $M = (I - W)^{T}(I - W)$

7. 取 a_1, a_2, \cdots, a_d 为对应的 d 个最小特征值对应的特征向量,得到线性变换 矩阵 α

3 仿真与实验

为评估提出算法的有效性,将 NFS-LPP 和 NFS-NPE 算法分别应用到 Yale 和 ORL 人脸库中进行分类和聚类实验仿真。在 LPP 和 NPE 算法中,只存在参数邻域大小 k 的选择,实验中邻域大小选择为 5。在 NFS-LPP 和 NFS-NPE 算法中,随机产生不同数量的正约束对,利用式(4)选择一定比例特征构成特征子空间(在 Yale 人脸库中,选择保留 20% 的特征;在 ORL 人脸库中,选择保留 40% 的特征)。在 Yale 和 ORL 人脸库中,分别随机选择每人6 张和5 张人脸图像作为训练样本,剩余人脸图像作为测试样本。分类器采用 KNN 比较分类正确率(为减小随机正约束产生的影响,分类正确率为 100 次实验的平均结果)。

在 Yale 人脸库中,LPP^[7] 和 NPE^[8] 算法的分类正确率分别为 74.43% 和 75.41%;在 ORL 人脸库中,其分类正确率分别为 84.94% 和 88.00%。表 2 中给出了在不同正约束数量下 NFS-LPP 和 NFS-NPE 的分类性能,并与 SSML^[9] 和 IDSDRC^[10] 算法进行了比较。SSML 和 IDSDRC 是近年提出的两种结合成对约束信息和无标签数据的半监督降维方法。SSML 在保持局部邻域结构的同时,进一步拉大负约束关系的距离和缩小正约束关系对间的距离,并通过参数 α 来控制局部邻域结构的贡献度;IDSDRC 在 SSML 的基础上进一步利用成对约束信息构建约束块,使得块类样本紧聚,块间样本散开,并增加正负约束距离控制参数 β 和 γ 。由于本文算法仅仅利用了正约束信息,而在 SSML 和 IDSDRC 中需要同时结合正负约束信息,为

了在同等正约束下进行算法的性能比较,仿真时同时为 SSML 和 IDSDRC 额外产生了相应比例的负约束信息。

表2 NFS-LPP 和 NFS-NPE 在不同正约束下的性能 (a) Yale 人脸库(6trains, 正确率/%)

| 算法 | 正约束对数 | | | | | | | | | |
|----------|-------|-------|-------|-------|-------|-------|-------|--|--|--|
| | 20 | 30 | 40 | 50 | 60 | 80 | 100 | | | |
| NFS-LPP_ | 82.42 | 82.27 | 82.83 | 82.96 | 83.40 | 83.47 | 83.92 | | | |
| NFS-LPP | 83.00 | 82.86 | 83.79 | 84.23 | 84.95 | 85.74 | 86.33 | | | |
| NFS-NPE_ | 79.03 | 78.93 | 79.51 | 79.48 | 79.99 | 79.88 | 80.03 | | | |
| NFS-NPE | 79.56 | 79.76 | 81.43 | 82.68 | 84.10 | 85.72 | 88.20 | | | |
| SSML | 56.68 | 60.60 | 65.87 | 71.12 | 73.96 | 85.26 | 90.11 | | | |
| IDSDRC | 70.80 | 63.95 | 67.88 | 75.74 | 82.44 | 88.64 | 90.75 | | | |

(b) ORL 人脸库(5 trains,正确率/%)

| 算法 | 正约束对数 | | | | | | | | |
|----------|-------|-------|-------|--------|-------|-------|-------|--|--|
| | 20 | 30 | 40 | 50 | 60 | 80 | 100 | | |
| NFS-LPP_ | 88.96 | 88.92 | 89.03 | 89.09 | 89.18 | 89.46 | 89.56 | | |
| NFS-LPP | 89.13 | 89.54 | 89.60 | 90.10 | 90.26 | 91.09 | 92.32 | | |
| NFS-NPE_ | 89.56 | 89.88 | 90.13 | 90.17 | 90.16 | 90.27 | 90.44 | | |
| NFS-NPE | 90.12 | 90.13 | 90.42 | 90.53 | 90.44 | 91.27 | 91.95 | | |
| SSML | 43.29 | 48.97 | 49.88 | 60.85 | 71.04 | 87.24 | 92.40 | | |
| IDSDRC | 71.52 | 68.86 | 67.67 | 74. 17 | 81.04 | 91.37 | 94.47 | | |

NFS-LPP_和 NFS-NPE_算法是利用正约束选择的特征子 空间邻域,与 LPP 和 NPE 算法的唯一区别就是在不同的空间 来寻找邻域结构。从表 2 中可以看出, 当随机产生 20 对正约 東时,NFS-LPP_在 Yale 和 ORL 人脸库中的分类正确率与 LPP 相比,分别提高了8%和4%左右,而NFS-NPE_的分类正确率 与 NPE 相比,也分别提高了 4% 和 2% 左右。这充分说明了在 特征子空间中构建邻域关系提高了数据流形的可区分性。但 表 2 同时表明, 当正约束对增加时, NFS-LPP_和 NFS-NPE_的 正确率增加并不明显,而 NFS-LPP 和 NFS-NPE 由于进一步利 用了正约束修正邻域矩阵权值和邻域成员,使得分类正确率随 之增加。这也验证了2.2节所述,当正约束的数量达到一定程 度后,不同约束数量下选择出来的特征趋于一致,以至圩构建 的邻域结构也趋于稳定。当可利用正约束较多时,单纯利用正 约束选择特征子空间邻域并不能充分利用约束信息。从表 2 还可以看出,NFS-LPP 和 NFS-NPE 在成对约束信息较少时,取 得了较高的分类正确率,远远优于 SSML 和 IDSDRC 算法。只 有在成对约束信息达到一定数量时,SSML 和 IDSDRC 分类性 能才略优于 NFS-LPP 和 NFS-NPE。

为进一步说明本文算法的有效性,分别在 Yale 和 ORL 人脸库中进行了 K-means 聚类实验。图 2 给出了各算法在不同正约束数量下的聚类性能比较。从图中可以看出,当可利用的成对约束信息较少时,SSML 和 IDSDRC 的聚类性能远低于无监督学习的算法,只有当约束达到一定数量时,才能获得较好的聚类性能。而 NFS-LPP 和 NFS-NPE 算法能一直保持较好的聚类性能,且随着正约束对的增加,聚类性能也随之提升。

结合分类和聚类实验可以看出:在基于正约束选择的特征子空间中构造邻域非常有效,尤其是它能利用少量约束信息选择出一个合适的特征子空间来表征内在的邻域结构。当约束较多时,NFS-LPP和NFS-NPE能充分利用约束信息来提升算法的性能。与SSML和IDSDRC算法的比较更体现出了本文算法的优越性,SSML和IDSDRC的性能只有当约束较多时,才能和本文算法相抗衡,而现实中往往可利用的约束信息是较少的。另外,在SSML和IDSDRC算法中存在多参数选择的问题,且参数选择无理论性指导,而本文提出的NFS-LPP和NFS-

NPE 算法只有原 LPP 和 NPE 算法中的参数邻域大小一个参数,并没有增加额外参数。

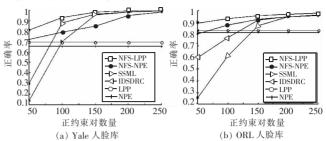


图 2 不同算法在不同数量正约束下的聚类性能比较

4 结束语

本文针对于高维数据的稀疏性和可能存在的非线性结构,以及一些高维数据难以获取足够多的训练样本来估计内在流形的问题,提出了一种利用已知先验正约束信息选择特征子空间的方法,并依据特征子空间的邻域结构来挖掘高维数据的内在流形,构建了基于特征子空间邻域特性的局部保持流形学习算法(NFS-LPP 和 NFS-NPE)。在 Yale 和 ORL 人脸库中的分类和聚类实验及与其他半监督流形学习算法对比验证了算法的有效性和优越性。但是,对于选择多少比例的特征来构造特征子空间,本文没有给出理论性指导,这也是未来有待进一步研究的问题。

参考文献:

- [1] TURK M, PENTLAND A. Eigenfaces for recognition [J]. Journal of Cognitive Neuroscience, 1991, 3(1);71-86.
- [2] BELHUMEUR P N, HEPANHA J P, KRIEGMAN D J. Eigenfaces vs. fisherfaces: recognition using class specific linear projection [J]. IEEE Trans on Pattern Analysis and Machine Intelligence, 1997, 19 (7):711-720.
- [3] TENENBAUM J B , De SILVA V, LANGFORD J C. A global geometric framework for nonlinear dimensionality reduction [J]. Science, 2000, 290 (5500); 2319-2323.
- [4] SHAO Ji-dong, RONG Gang. Nonlinear process monitoring based on maximum variance unfolding projections [J]. Expert Systems with Applications, 2009, 36(8):11332-11340.
- [5] PAN Yao-zhang, GE S S, MAMUNA A A. Weighted locally linear embedding for dimension reduction [J]. Pattern Recognition, 2009, 42 (5):798-811.
- [6] BELKIN M, NIYOGI P. Laplacian eigenmaps for dimensionality reduction and data representation [J]. Neural Computation, 2003, 15 (6):1373-1396.
- [7] HE Xiao-fei, NIYOGI P. Locality preserving projections [C]//Proc of Advances in Neural Information Processing Systems. Cambridge, MA: MIT Press, 2004:153-160.
- [8] HE Xiao-fei, CAI Deng, YAN Shui-cheng, et al. Neighborhood preserving embedding [C]//Proc of the 10th International Conference on Computer Vision. Washington DC: IEEE Computer Society, 2005: 1208-1213.
- [9] BAGHSHAH M S, SHOURAKI B S. Semi-supervised metric learning using pairwise constraints [C]//Proc of the 21st International Joint Conference on Artifical Intelligence. San Francisco, CA: Morgan Kaufmann Publishers Inc, 2009:1217-1222.
- [10] WANG Na, LI Xia, CUI Ying-jie, *et al.* Instance-level based discriminative semi-supervised dimensionality reduction with chunklets [J]. International Journal of Innovative Computing, Information and Control, 2010, 6(8):3763-3773.