

文章编号: 1001-0920(2011)10-1511-04

## 基于 greedy GDA 的训练数据减少和非线性特征提取方法

刘小芳<sup>1,2</sup>, 何彬彬<sup>1</sup>, 李小文<sup>1</sup>

(1. 电子科技大学 地表空间信息技术研究所, 成都 611731; 2. 四川理工学院 计算机学院, 四川 自贡 643000)

**摘要:** 针对大数据集情况下标准广义判别分析(GDA)方法进行非线性特征提取时计算复杂度较高的问题, 提出了基于 GGDA (greedy GDA) 的训练数据减少和非线性特征提取方法. 该方法用 greedy 核主成分分析方法的 greedy 技术对训练数据选取子集; 然后用 GDA 方法对子集而不是全部训练数据训练特征提取模型; 并用几种特征提取的数据进行分类对比实验. 实验结果表明, GGDA 和 GDA 方法的特征提取性能优于其他对比方法, GGDA 方法不仅较好地保持了 GDA 方法的特征提取性能, 而且减少了大数据集进行非线性特征提取的计算复杂度.

**关键词:** greedy 广义判别分析; greedy 核主成分分析; 训练数据减少; 非线性特征提取; 核矩阵; 分类  
**中图分类号:** TP391      **文献标识码:** A

## Greedy GDA method for training data reduction and nonlinear feature extraction

LIU Xiao-fang<sup>1,2</sup>, HE Bin-bin<sup>1</sup>, LI Xiao-wen<sup>1</sup>

(1. Institute of Geo-spatial Information Science and Technology, University of Electronic Science and Technology of China, Chengdu 611731, China; 2. Department of Computer Science and Technology, Sichuan University of Science and Engineering, Zigong 643000, China. Correspondent: LIU Xiao-fang, E-mail: lxf1969@163.com)

**Abstract:** Nonlinear feature extraction using standard generalized discriminant analysis(GDA) has high computational complexity in large datasets. Therefore, a greedy GDA(GGDA) is proposed to reduce training data and deal with the nonlinear feature extraction problem. Firstly, a subset is selected from the full training data by using the greedy technique of the greedy KPCA(GKPCA) method. Then, the feature extraction model is trained by using the GDA method with the subset instead of the full training data. Finally, classification experiments using data of several feature extraction methods are performed. The simulation results show that the feature extraction performance of both the GGDA and the GDA methods outperform that of other methods. In addition of retaining the performance of the GDA method, the GGDA method reduces the computational complexity of the nonlinear feature extraction in large datasets.

**Key words:** greedy generalized discriminant analysis; greedy kernel principal component analysis; training data reduction; nonlinear feature extraction; kernel matrix; classification

### 1 引言

近年来, 核理论成为机器学习和模式识别的一个快速发展分支. 核理论利用满足 Mercer 条件的核函数巧妙地推导出线性算法的非线性形式<sup>[1]</sup>, 如: 核主成分分析(KPCA)<sup>[2]</sup>和广义判别分析(GDA)<sup>[3]</sup>是通过核理论分别对主成分分析(PCA)和线性判别分析(LDA)方法进行非线性扩展而形成的, 显示了较强的非线性特征提取能力.

GDA 和 KPCA 方法应用核理论能够较好地提取非线性特征. 但两种方法进行特征提取时均需要计算

和存储核矩阵, 并对核矩阵进行特征值分解. 核矩阵的维数等于样本数的平方, 核矩阵所需的计算时间和内存与输入空间的维数基本无关, 而与样本数目密切相关, 随着样本数目的增加, 内存与计算量存在问题, 使问题的求解更加困难, 甚至无法实现<sup>[4-5]</sup>. 因此, 有必要减少大数据集训练数据. Franc 等人<sup>[6]</sup>在 GKPCA (greedy KPCA) 方法中提出用 greedy 技术寻找一个低维的数据表示近似再生核 Hilbert 空间的大数据集. 类似于 KPCA 方法, GKPCA 方法的目的是最小化近似数据的均方重构误差; 不同之处在于,

收稿日期: 2010-06-01; 修回日期: 2010-10-07.

基金项目: 国家自然科学基金项目(40701146); 国家 863 计划项目(2007AA12Z227).

作者简介: 刘小芳(1969—), 女, 副教授, 博士生, 从事智能信息处理、模式识别等研究; 李小文(1947—), 男, 中国科学院院士, 教授, 博士生导师, 从事定量遥感反演等研究.

GKPCA 方法用于表示数据的低维空间的基向量是从训练数据集中选取的向量,而不是全部训练数据的线性组合,并找出了描述原始数据的简单模型.对比 KPCA 方法, GKPCA 方法根据重构误差找到了数据表示的次优解,该方法可用于降低核方法的计算量和存储需求,同时也可以用于降低用核方法进行学习的复杂度,还能进行在线特征提取.由于 GKPCA 方法的计算以内积的形式出现,可以引入核方法解决非线性问题.

KPCA 和 GDA 方法均是将输入空间的数据映射到具有非线性投影方向的再生核 Hilbert 特征空间. KPCA 方法是提取互不相关的特征向量,对于分类问题而言不一定是最佳的.而 GDA 方法则考虑了类内与类间的信息,这种特征提取方法对于分类问题而言是较好的<sup>[7]</sup>.因此,本文将 GKPCA 方法中提取子集的 greedy 技术与 GDA 方法相结合,形成 GGDA (Greedy GDA) 方法,用于训练数据的减少和非线性特征提取. GGDA 方法是在使用 GDA 方法之前,先用 GKPCA 方法中提取子集的 greedy 技术对训练数据选取子集,子集能较好地表示原始数据,用子集而不是全部训练数据训练特征提取模型.

## 2 Greedy GDA 方法

### 2.1 greedy 技术对训练数据的减少

GKPCA 方法是标准 KPCA 的一种有效实现.给定训练数据集  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \subset \mathbf{R}^q$ , GKPCA 方法的目的是从训练数据集  $\mathbf{X}$  中找到一个子集  $\mathbf{X}_S = \{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_n\} \subset \mathbf{R}^q$  (即  $\mathbf{X}_S \subseteq \mathbf{X}$ ), 使  $\text{span}(\mathbf{X}_S) \approx \text{span}(\mathbf{X})$ .  $\mathbf{I} = \{i_1, i_2, \dots, i_N\}$  是训练数据集  $\mathbf{X}$  的索引矩阵,  $\mathbf{J} = \{j_1, j_2, \dots, j_n\}$  是子集  $\mathbf{X}_S$  的索引矩阵,且  $\mathbf{J} \subseteq \mathbf{I}$ .

GKPCA 方法是在再生核 Hilbert 空间实现,训练数据集  $\mathbf{X}$  通过一个非线性映射  $\Phi$  变换到特征空间  $\mathbf{F}$ .  $\{\Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2), \dots, \Phi(\mathbf{x}_N)\}$  和  $\{\Phi(\mathbf{s}_1), \Phi(\mathbf{s}_2), \dots, \Phi(\mathbf{s}_n)\}$  分别是训练数据集  $\mathbf{X}$  和子集  $\mathbf{X}_S$  在特征空间的非线性映射. GKPCA 方法是在特征空间  $\mathbf{F}$  中找到一种新的核扩展,能较好地近似原始的核扩展,即  $\text{span}(\Phi(\mathbf{X}_S)) \approx \text{span}(\Phi(\mathbf{X}))$ . 原始训练数据  $\mathbf{X}$  在近似特征空间表示为

$$\tilde{\Phi}(\mathbf{x}_i) = \sum_{j \in \mathbf{J}} [\beta_i]_j \Phi(\mathbf{x}_j), \quad i \in \mathbf{I}, \quad (1)$$

其中  $\beta_i \in \mathbf{R}^n$  是一组系数. 训练数据集  $\mathbf{X}$  在特征空间的近似表示  $\{\tilde{\Phi}(\mathbf{x}_1), \tilde{\Phi}(\mathbf{x}_2), \dots, \tilde{\Phi}(\mathbf{x}_N)\}$  为原始数据空间子集  $\mathbf{X}_S$  在特征空间  $\{\Phi(\mathbf{s}_1), \Phi(\mathbf{s}_2), \dots, \Phi(\mathbf{s}_n)\}$  内的线性组合.

寻找近似的核扩展问题可以表示为最优化问题.

当子集  $\mathbf{X}_S$  的样本数  $n$  很小时,最小化均方重构误差

$$\varepsilon_{\text{MS}}(\mathbf{X}|\mathbf{J}) = \frac{1}{N} \sum_{i \in \mathbf{I}} \left\| \Phi(\mathbf{x}_i) - \sum_{j \in \mathbf{J}} [\beta_i]_j \Phi(\mathbf{x}_j) \right\|^2. \quad (2)$$

训练数据集  $\mathbf{X}$  的重构误差  $\varepsilon_{\text{MS}}(\mathbf{X}|\mathbf{J})$  依赖于子集  $\mathbf{X}_S$  和系数  $\beta_i \in \mathbf{R}^n (i \in \mathbf{I})$ . 给定子集  $\mathbf{X}_S$ , 系数  $\beta_i$  通过最小化均方重构误差  $\varepsilon_{\text{MS}}(\mathbf{X}|\mathbf{J})$  得到,即

$$\beta_i = \arg \min_{\beta \in \mathbf{R}^n} \varepsilon_{\text{MS}}(\mathbf{X}|\mathbf{J}), \quad i \in \mathbf{I}. \quad (3)$$

为了避免训练数据集  $\mathbf{X}$  通过非线性映射函数  $\Phi(\mathbf{x}_i)$  显式地映射到特征空间  $\mathbf{F}$ , 用核函数重新表示式(2)为

$$\varepsilon_{\text{MS}}(\mathbf{X}|\mathbf{J}) = \frac{1}{N} \sum_{i \in \mathbf{I}} (k(\mathbf{x}_i, \mathbf{x}_i) - 2\mathbf{K}_S \mathbf{k}_S(\mathbf{x}_i) + \langle \mathbf{k}_S(\mathbf{x}_i), \mathbf{K}_S \mathbf{k}_S(\mathbf{x}_i) \rangle).$$

其中:  $\mathbf{K}_S \in \mathbf{R}^{n \times n}$  是子集  $\mathbf{X}_S$  的核矩阵;  $\mathbf{k}_S(\mathbf{x}_i) = [k(\mathbf{x}_{j_1}, \mathbf{x}_i), k(\mathbf{x}_{j_2}, \mathbf{x}_i), \dots, k(\mathbf{x}_{j_n}, \mathbf{x}_i)]^T \in \mathbf{R}^n$  是样本  $\mathbf{x}_i$  映射到子集  $\mathbf{X}_S$  的特征空间  $\mathbf{F}$  的投影.

GKPCA 方法需要从训练数据集  $\mathbf{X}$  寻找最佳的子集  $\mathbf{X}_S$ , 即从  $\mathbf{X}$  的索引矩阵  $\mathbf{I}$  中找出子集  $\mathbf{X}_S$  的索引矩阵  $\mathbf{J}$ , 有

$$\mathbf{J} = \arg \min_{\mathbf{J} \subseteq \mathbf{I}} \varepsilon_{\text{MS}}(\mathbf{X}|\mathbf{J}), \quad (4)$$

且  $\text{card}(\mathbf{J}) = n$ . 如果给定了子集  $\mathbf{X}_S$  的样本数  $n$ , 通过最小化式(4)求解子集  $\mathbf{X}_S$  的问题有  $\binom{N}{n}$  组合, 这是不切实际的. 文献[8]提出求解最小化  $\varepsilon_{\text{MS}}(\mathbf{X}|\mathbf{J})$  上界, 即

$$\varepsilon_{\text{MS}}(\mathbf{X}|\mathbf{J}) \leq \frac{1}{N} (N - n) \sum_{i \in \mathbf{I} \setminus \mathbf{J}} \max \|\Phi(\mathbf{x}_i) - \tilde{\Phi}(\mathbf{x}_i)\|^2, \quad (5)$$

其中  $(N - n)$  为尺度系数, 是由于子集  $\mathbf{X}_S$  里的  $n$  个样本有零表示误差. 式(5)一定有效, 因为均方重构误差  $\varepsilon_{\text{MS}}(\mathbf{X}|\mathbf{J})$  不高于最大均值误差. 通过 greedy 技术实现式(5), 具体实现过程见文献[8], 每次迭代选择一个最大误差的样本进入子集  $\mathbf{X}_S$ .

### 2.2 GDA 方法的非线性特征提取

GDA 方法是将原始训练数据通过非线性映射  $\Phi$  变换到特征空间  $\mathbf{F}$  中, 然后在特征空间  $\mathbf{F}$  中完成 LDA 变换. 这里, GDA 方法的原始训练数据为通过 GKPCA 方法中的 greedy 技术从训练数据集  $\mathbf{X}$  中选取的子集  $\mathbf{X}_S$ . 为了便于理解, 将子集  $\mathbf{X}_S$  重新定义为  $\mathbf{T}_{\mathbf{X}\mathbf{Y}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$ ,  $\mathbf{x}_i \in \mathbf{R}^q$ ,  $y \in \mathbf{Y} = \{1, 2, \dots, c\}$ . GDA 方法的训练数据  $\mathbf{T}_{\mathbf{X}\mathbf{Y}}$  通过非线性映射函数  $\Phi: \mathbf{X} \rightarrow \mathbf{F}$ , 在特征空间  $\mathbf{F}$  中对应  $\mathbf{T}_{\Phi\mathbf{Y}} = \{(\Phi(\mathbf{x}_1), y_1), (\Phi(\mathbf{x}_2), y_2), \dots, (\Phi(\mathbf{x}_n), y_n)\}$ . 训练数据集  $\mathbf{T}_{\Phi\mathbf{Y}}$  在特征空间  $\mathbf{F}$  中的类内离散矩阵  $\mathbf{S}_w$  和类间离散矩阵  $\mathbf{S}_b$  分别为

$$\mathbf{S}_w = \frac{1}{n} \sum_{i=1}^c \sum_{j=1}^{n_i} (\Phi(\mathbf{x}_{ij}) - \boldsymbol{\mu}_i)(\Phi(\mathbf{x}_{ij}) - \boldsymbol{\mu}_i)^T, \quad (6)$$

$$\mathbf{S}_b = \frac{1}{n} \sum_{i=1}^c \frac{n_i}{n} (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T. \quad (7)$$

其中:  $n_i$  为第  $i$  类训练样本的个数,  $\sum_{i=1}^c n_i = n$ ;  $\Phi(\mathbf{x}_{ij})$  为特征空间  $\mathbf{F}$  中第  $i$  类第  $j$  个训练样本;  $\Phi(\mathbf{x}_j)$  为特征空间  $\mathbf{F}$  中第  $j$  个训练样本;  $\boldsymbol{\mu}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} \Phi(\mathbf{x}_j)$  为特征空间  $\mathbf{F}$  中第  $i$  类训练样本的均值;  $\boldsymbol{\mu} = \frac{1}{n} \sum_{j=1}^n \Phi(\mathbf{x}_j)$  为特征空间  $\mathbf{F}$  中全部训练样本的均值。

GDA 的目的在于找到一组向量  $\mathbf{u} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n]$ , 使类间离散矩阵  $\mathbf{S}_b$  和类内离散矩阵  $\mathbf{S}_w$  之比达到最大, 即

$$\mathbf{u} = \arg \max_{\mathbf{u}} \frac{|\mathbf{u}^T \mathbf{S}_b \mathbf{u}|}{|\mathbf{u}^T \mathbf{S}_w \mathbf{u}|}. \quad (8)$$

式 (8) 的解实际是找到矩阵  $\mathbf{S}_w^{-1} \mathbf{S}_b$  的特征向量  $\mathbf{u}$  对应的最大特征值  $\lambda$ <sup>[9]</sup>, 即

$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{u} = \lambda \mathbf{u}. \quad (9)$$

最大化商  $J(\mathbf{u}) = |\mathbf{u}^T \mathbf{S}_b \mathbf{u}| / |\mathbf{u}^T \mathbf{S}_w \mathbf{u}|$ , 即对  $\mathbf{u}$  求导数, 有

$$\begin{aligned} \frac{\partial J(\mathbf{u})}{\partial \mathbf{u}} = 0 &\Rightarrow \mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{u} = \left( \frac{\mathbf{u}^T \mathbf{S}_b \mathbf{u}}{\mathbf{u}^T \mathbf{S}_w \mathbf{u}} \right) \mathbf{u} \Rightarrow \\ \lambda &= \left( \frac{\mathbf{u}^T \mathbf{S}_b \mathbf{u}}{\mathbf{u}^T \mathbf{S}_w \mathbf{u}} \right). \end{aligned} \quad (10)$$

特征向量  $\mathbf{u} \in \mathbf{F}$  在  $\Phi(\mathbf{x}_{ij})$  张成的空间里, 存在系数  $\boldsymbol{\alpha} = (\alpha_{ij})_{i=1,2,\dots,c, j=1,2,\dots,n_i}$ , 使得

$$\mathbf{u} = \sum_{i=1}^c \sum_{j=1}^{n_i} \alpha_{ij} \Phi(\mathbf{x}_{ij}). \quad (11)$$

求解特征向量  $\mathbf{u}$  可转化为求相应的系数  $\boldsymbol{\alpha}$ . 为了能够应用核函数  $k(\mathbf{x}_i, \mathbf{x}_j) = k_{ij} = \Phi^T(\mathbf{x}_i) \Phi(\mathbf{x}_j)$  使训练数据集  $\mathbf{X}$  通过非线性映射函数  $\Phi$  隐式地映射到特征空间, 用  $\Phi^T(\mathbf{x}_{ij})$  乘以式 (9), 得到

$$\Phi^T(\mathbf{x}_{ij}) \mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{u} = \lambda \Phi^T(\mathbf{x}_{ij}) \mathbf{u}. \quad (12)$$

式 (12) 和 (9) 有相同的特征向量<sup>[2]</sup>. 将式 (11) 代入 (12) 得

$$\lambda = \left( \frac{\boldsymbol{\alpha}^T \mathbf{K} \mathbf{W} \mathbf{K} \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T \mathbf{K} \mathbf{K} \boldsymbol{\alpha}} \right). \quad (13)$$

其中

$$\mathbf{K} = (\mathbf{K}_{pq})_{p=1,2,\dots,c, q=1,2,\dots,c}, \quad (14)$$

$$\mathbf{W} = (\mathbf{W}_i)_{i=1,2,\dots,c}. \quad (15)$$

核矩阵  $\mathbf{K}$  是一个分块矩阵, 其子矩阵  $\mathbf{K}_{pq} = (k_{ij})_{p=1,2,\dots,n_p, q=1,2,\dots,n_q}$  由第  $p$  类和第  $q$  类的数据经过内积  $(k_{ij})_{pq} = \Phi^T(\mathbf{x}_{pi}) \Phi(\mathbf{x}_{qj})$  得到,  $\mathbf{K}_{pq}$  是  $n_p \times n_q$  矩阵,  $\mathbf{K}$  是  $n \times n$  的对称矩阵. 系数矩阵  $\mathbf{W}$  是一个

$n \times n$  的分块对角阵, 其子矩阵  $\mathbf{W}_i$  是所有元素为  $1/n_i$  的  $n_i \times n_i$  矩阵.

对核矩阵  $\mathbf{K}$  进行特征值分解, 求出  $\boldsymbol{\alpha}$ , 具体求解过程见文献 [3]. 通过归一化  $\boldsymbol{\alpha}$  使特征向量  $\mathbf{u}$  归一化, 即将式 (11) 代入  $\mathbf{u}^T \mathbf{u} = 1$ , 得到

$$\mathbf{u}^T \mathbf{u} = \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} = 1. \quad (16)$$

系数  $\boldsymbol{\alpha}$  除以  $\sqrt{\boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha}}$ , 得到归一化的特征向量  $\mathbf{u}$ . 测试数据  $\mathbf{z}$  在特征向量  $\mathbf{u}$  上的投影为

$$\begin{aligned} \mathbf{u}^T \Phi(\mathbf{z}) &= \sum_{i=1}^c \sum_{j=1}^{n_i} \alpha_{ij} \Phi(\mathbf{x}_{ij}) \Phi(\mathbf{z}) = \\ &= \sum_{i=1}^c \sum_{j=1}^{n_i} \alpha_{ij} k(\mathbf{x}_{ij}, \mathbf{z}). \end{aligned} \quad (17)$$

算法要求非线性映射函数  $\Phi(\mathbf{x}_j)$  中心化, 由于  $\Phi(\mathbf{x}_j)$  是隐式映射到特征空间  $\mathbf{F}$ , 其中中心化是通过中心化特征空间  $\mathbf{F}$  的核矩阵  $\mathbf{K}$  实现的<sup>[3]</sup>.

### 2.3 GGDA 方法的实现步骤

GGDA 方法的训练数据减少和非线性特征提取的实现步骤如下:

- 1) 通过 GKPCA 方法中的 greedy 技术实现式 (5), 从训练数据集  $\mathbf{X}$  中选取子集  $\mathbf{X}_S$  作为 GDA 方法的训练数据  $\mathbf{T}_{\mathbf{X}\mathbf{Y}}$ ;
- 2) 用式 (14) 和 (15) 分别计算核矩阵  $\mathbf{K}$  和系数矩阵  $\mathbf{W}$ ;
- 3) 在特征空间对核矩阵  $\mathbf{K}$  进行特征值分解, 计算其特征值并标准化特征向量;
- 4) 找出最大特征值  $\lambda$  和对应的特征向量  $\mathbf{u}$ ;
- 5) 根据式 (17) 计算测试样本  $\mathbf{z}$  在特征向量  $\mathbf{u}$  上的投影.

对比 GDA 方法中核矩阵  $\mathbf{K}$  的大小 ( $N \times N$ ), GGDA 方法核矩阵  $\mathbf{K}$  的大小为  $(n \times n)$ ; GDA 方法和 GGDA 方法的时间复杂度分别为  $O(N^3)$  和  $O(Nn^2)$ ; 当  $n \ll N$  时, GGDA 方法大大降低了 GDA 方法的内存需求和计算时间.

## 3 实验结果及分析

实验平台为 Intel 酷睿 2 双核 T6500 2.10 GHz 处理器, 2 G 内存, Windows XP 操作系统. 实验算法用 C 和 Matlab 混合编程实现. 实验数据选取 UCI 数据库数据集 Iris 和 pendigits<sup>[10]</sup>. Iris 数据个数较少, pendigits 数据个数较多. 对于数据集 Iris, 数据特点为一类数据与其他类数据离得较远, 另外两类数据离得较近, 而且有部分数据交叠. 数据集 pendigits 为 0~9 数字的钢笔字识别数据库. 表 1 给出了数据集的描述.

为了验证 GGDA 方法的性能, 对于几种特征提取方法提取的特征数据, 用  $k$ -最近邻分类算法<sup>[11]</sup>进

表 1 数据集的描述

数据集	样本数	类别	属性
Iris	150	3	4
pendigits	10992	9	16

行分类对比分析.  $k$ -近邻分类器(kNN)是从测试样本  $\mathbf{x}_i$  开始生长, 不断地扩大区域, 直到包含进  $k$  个距离最近的训练样本为止, 并将测试样本  $\mathbf{x}_i$  的类别归为最近的  $k$  个训练样本中出现频率最大的类别. 这里用欧氏距离作为  $k$ -近邻分类器设计中衡量样本之间距离的度量函数.

对于给定的训练数据集  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \subset \mathbf{R}^q$ , 通过 GKPCA 方法的 greedy 技术从训练数据集  $\mathbf{X}$  中找到子集  $\mathbf{X}_S = \{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_n\} \subset \mathbf{R}^q$ . LDA 和 PCA 方法从  $\mathbf{X}$  中最多提取  $q$  个特征; GDA 和 KPCA 方法从  $\mathbf{X}$  中最多提取  $N$  个特征; GGDA 和 GKPCA 方法从  $\mathbf{X}_S$  中最多提取  $n$  个特征, 有  $q < n \leq N$ . 本文主要目的是验证子集  $\mathbf{X}_S$  的非线特征提取性能, 没有考虑不同特征个数对分类的影响, 因此, 分类实验统一选取变换后的  $q$  个特征.

GGDA, GDA, GKPCA 和 KPCA 方法的核函数取高斯核函数  $\mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / (2\sigma^2))$ . 通过 5-重交叉验证法<sup>[12]</sup>, 对于数据集 Iris 和 pendigits,  $k$ -最近邻分类算法的  $k$  分别取 9 和 3; 核参数  $\sigma$  分别取 4 和 60; 用 GKPCA 方法提取的子集  $\mathbf{X}_S$  的样本数  $n$  分别为 17 和 200. 分类结果用训练分类器的时间  $T_c$ , 训练数据识别率以及测试数据的识别率进行了评价. 对于 Iris 和 pendigits 数据, 不同特征提取方法得到的特征提取数据用  $k$ -近邻分类算法分类的结果分别

表 2 几种方法对 Iris 数据的分类性能比较

方法	$T_c(s)$	识别率/%	
		训练	测试
kNN	0.009 050	93.33	90.00
PCA+kNN	0.009 245	95.37	90.16
KPCA+kNN	0.009 494	95.36	90.33
GKPCA+kNN	0.000 696	95.33	90.26
LDA+kNN	0.001 237	96.67	96.33
GDA+kNN	0.001 542	98.78	96.67
GGDA+kNN	0.000 323	98.73	96.59

表 3 几种方法对 pendigits 数据的分类性能比较

方法	$T_c(s)$	识别率/%	
		训练	测试
kNN	5.49	98.17	94.97
PCA+kNN	5.53	98.35	95.17
KPCA+kNN	5.87	99.35	96.57
GKPCA+kNN	0.59	99.18	96.29
LDA+kNN	3.12	98.71	95.45
GDA+kNN	3.18	99.77	97.83
GGDA+kNN	0.38	99.58	97.71

如表 2 和表 3 所示, 表中的评价结果是利用 5-重交叉验证法得到的平均值.

2 个数据集和 6 种特征提取数据的分类结果表明: 用特征提取的数据进行分类比不用特征提取的数据进行分类的性能要好; 用线性特征提取方法 LDA 特征提取的数据分类结果好于 PCA 方法; 用非线性特征提取方法 GGDA 和 GDA 的特征提取性能优于 GKPCA 和 KPCA 方法; 核空间的非线性特征提取方法 GGDA, GDA, GKPCA 和 KPCA 优于输入空间的线性特征提取方法 LDA 和 PCA; GKPCA 方法的性能接近于 KPCA 方法的性能; GGDA 方法的性能接近于 GDA 方法的性能. 总之, 在实验中不管对于大数据集还是小数据集, GGDA 和 GDA 方法的特征提取性能均优于其他方法. GGDA 方法不仅较好地保持了 GDA 方法特征提取性能, 而且减少了大数据集进行非线性特征提取的内存和计算时间. 实验也表明, 由于 GGDA 和 GDA 方法考虑了类内与类间的信息, 这种特征提取方法对于分类问题而言是较好的. 因此, 尽管 GGDA 方法是根据重构误差找到数据表示的次优解, 但该方法可用于降低核方法的计算量和存储需求, 同时也可以用于降低用核方法进行学习的复杂度, 还能进行在线非线性特征的提取.

## 4 结 论

本文提出将一种有效的特征提取方法 GGDA 方法用于训练数据减少和非线性特征的提取, 该方法有效地结合了 GKPCA 方法中提取子集的 greedy 技术和非线性特征提取 GDA 方法. 对比实验结果表明, GGDA 和 GDA 方法的特征提取性能优于其他方法. GGDA 方法不仅较好地保持了 GDA 方法的非线性特征提取性能, 而且减少了大数据集进行非线性特征提取的计算复杂度. 由于 GGDA 和 GDA 方法考虑了类内与类间的信息, 这种特征提取方法对于分类问题而言好于提取互不相关的特征向量的 GKPCA 和 KPCA 方法. 下一步研究方向是用 GGDA 方法对大数据集的在线非线性特征进行提取.

## 参考文献(References)

- [1] Schölkopf B, Smola A J. Learning with kernels[M]. Cambridge: MIT Press, 2002: 53-55.
- [2] Schölkopf B, Smola A, Muller K R. Nonlinear component analysis as a kernel eigenvalue problem[J]. Neural Computation, 1998, 10(5): 1299-1319.
- [3] Baudat G, Anouar F. Generalized discriminant analysis using a kernel approach[J]. Neural Computation, 2000, 12(10): 2385-2404.