基于广义判别分析的光谱分类

许 馨^{1,2},杨金福¹,吴福朝¹,赵永恒²

1. 中国科学院自动化研究所国家模式识别实验室, 北京 100080

2. 中国科学院国家天文台,北京 100012

摘 要 提出了基于广义判别分析(generalized discriminant analysis, GDA)方法对恒星(Star)、星系 (Galaxy)和类星体(Quasars)的光谱进行分类。广义判别分析将核技巧与 Fisher 判别分析结合起来,通过非 线性映射将样本集映射到高维特征空间 F,在 F 空间中进行线性判别分析。实验对比了 LDA,GDA,PCA, KPCA 算法对于恒星、星系和类星体的光谱分类性能。结果表明基于 GDA 的算法对于这 3 种类型光谱的分 类正确率最高,LDA 次之;尽管 KPCA 也是一种基于核的方法,但是选择主成分个数较少时效果较差,甚 至低于 LDA;基于 PCA 的分类效果最差。

主题词 光谱分类; 广义判别分析; 线性判别分析; 核主成分分析 中图分类号: TP29 文献标识码: A 文章编号: 1000 0593(2006)10196005

引 言

天文分类技术具有相对特殊的背景,目前使用的方法许 多是基于 PCA 的线性的方法[1,2] 和基于神经网络的非线性 方法[3]。近年来, 随着统计理论的发展, 基于核技术的学习 方法如 SVM^[4]、核主成分分析(KPCA)^[5] 和核 Fisher 判别分 析(kernel fisher discriminant, KFD)^[6]在模式识别与机器学 习领域的优越表现引起了各行业学者的广泛关注。因为核学 习方法既可以提取数据的非线性特征,又有较好的推广能 力。KPCA 由 Scholkopf^[5] 首先提出,随后 Mika^[6] 提出了基 于核的 Fisher 判别分析算法(KFD)。KPCA 是一种用于描述 数据的非线性的特征提取方法, KFD 是用于分类的非线性 的特征提取方法,后者更适合分类问题。KFD 的主要思想是 将样本数据集非线性映射到高维特征空间、以求数据能够线 性可分或近似线性可分,然后在核特征空间中进行线性 Fisher 判别分析, 找到最利于分类的投影方向, 从而导致输 入空间的非线性判别能力。Baudat 等^[7] 将 KFD 发展到多类 问题,称之为 generalized discriminant analysis(GDA)。

KFD(或者 GDA)能够抽取基于分类的非线性特征,它 不像神经网络那样依赖对模型的选择,而且也不存在维数灾 难和局部极小值问题,因而它在许多实际分类问题中非常有 效。在天文光谱的分类问题中,面对的数据依然是具有非线 性特征的,天体红移和类型的耦合使得我们在做分类时遇到 了很大的困难。由此,我们尝试用基于核技术的广义判别分析(GDA)对于恒星、星系和类星体进行分类,得到了较好的 实验结果。

本文结构如下,第1部分简单的介绍了多类的线性 Fisher 判别分析;第2部分的广义判别分析为基于核技术的 多类线性判别分析;第3部分是 KPCA 与 GDA 的联系;第4 部分为实验结果及分析;第5部分为全文总结。

1 多类的线性 Fisher 判别分析

Fisher 线性判别分析的本质在于找到一个子空间,使得 各个类别在这个子空间中能较好地分离,易于分类器的设 计^[8]。

假设有一组 $N \uparrow d$ 维样本的集合 $A = \{x_1, x_2, ..., x_N \in R^d\} = \bigcup_{j=1}^{c} A_j, d \ge c,$ 其中, c为类别数; A_j 是第j类的 $n_j \uparrow$ 样本构成的样本子集, $N = \sum_{j=1}^{c} n_j$ 。

常用的标准 Fisher 准则函数如下

$$J(W) = |W^{\mathrm{T}}S_{\mathrm{B}}W| / |W^{\mathrm{T}}S_{\mathrm{W}}W|$$
(1)

S_B, S_w分别称为类间散度矩阵和类内散度矩阵

$$S_{\mathrm{B}} = \sum_{k=1}^{c} n_{k} (\boldsymbol{\mu}_{k} - \boldsymbol{\mu}) (\boldsymbol{\mu}_{k} - \boldsymbol{\mu})^{\mathrm{T}}$$
$$S_{\mathrm{W}} = \sum_{k=1}^{c} \sum_{x_{i} \in A_{k}} (x_{i} - \boldsymbol{\mu}_{k}) (x_{i} - \boldsymbol{\mu}_{k})^{\mathrm{T}}$$

基金项目: 国家"863" 计划项目(2003 A A133060) 资助

作者简介: 许 馨, 女, 1974年生, 中国科学院自动化研究所国家模式识别实验室博士研究生

收稿日期: 2005-08-08, 修订日期: 2005-11-08

^{© 1994-2010} China Academic Journal Electronic Publishing House. All rights reserved. http://www.cnki.net

本的均值, $\mu = \frac{1}{N} \sum_{k=1}^{N} n_k \mu_k$ 。求解使得公式(1)最大化的矩阵 W 的列向量是下列等式中的非零特征值对应的特征向量 $S_{B}w_i = \lambda S_{W}w_i \Rightarrow (S_{B} - \lambda S_{W})w_i = 0$)

文献[9] 提出扩展的 Fisher 准则函数

 $J(W) = |W^{T}S_{B}W| / |W^{T}S_{T}W|$ (2) $S_{T} = S_{W} + S_{B}, S_{T}$ 称为总散度矩阵。在寻找最优投影方向时,公式(1)和(2)是等价的。

2 广义判别分析

广义判别分析(GDA)的思想是通过一个非线性映射,将 输入空间的样本映射到高维特征空间,在这个特征空间中进 行线性 Fisher 判决分析,从而实现相对于输入空间的非线性 判决分析^[6,7]。

$$\Phi$$
, $R^d \xrightarrow{\rightarrow} F$, $x \Phi | \xrightarrow{\rightarrow} (x)$

参照准则函数(2),在 F空间中进行线性 Fisher 判决, 其扩展的准则函数为

$$J(W^{\Phi}) = \arg \max_{W^{\Phi}} (|(W^{\Phi})^{\mathrm{T}} S^{\Phi}_{\mathrm{B}} W^{\Phi}|/|(W^{\Phi})^{\mathrm{T}} S^{\Phi}_{\mathrm{T}} W^{\Phi}|)$$

$$(3)$$

其中, $W^{\circ} \in F$; 假设数据在 F 空间已经中心化⁵, S_{Γ}° 和 S_{Γ}° 为相应的 F 空间的类间散度矩阵和总散度矩阵

$$S_{\rm B}^{\Phi} = \sum_{k=1}^{c} n_k (\mathcal{H}_k^{\Phi}) (\mathcal{H}_k^{\Phi})^{\rm T}$$
$$S_{\rm T}^{\Phi} = \frac{1}{N} \sum_{i=1}^{N} \Phi(x_i) \Phi(x_i)^{\rm T}$$
(4)

其中, $\mu_k^{\Phi} = \frac{1}{n_k} \sum_{l=1}^{n_k} \Phi(x_{kl})$ 。

求最优的 ₩[®] 等价于求解下述广义特征值和特征向量问 题:

$$\lambda_i S \stackrel{\text{def}}{=} W_i^{\Phi} = S_B \stackrel{\text{def}}{=} W_i^{\Phi}, \quad i = 1, \quad \dots, \quad m,$$

$$\lambda_1 \geq \lambda_2 \geq \dots \lambda_2 \geq \lambda_m, \ m \leq c-1 \tag{5}$$

注意到 W^{\oplus} 的列向量 W_i^{\oplus} 一定在 $\Phi(x_1)$, $\Phi(x_2)$, ..., $\Phi(x_N)$ 张 成的空间 *F* 中,因而存在系数 α_i (*j* = 1, ..., *N*),使得

$$W_i^{\Phi} = \sum_{j=1}^{N} \alpha_{ij} \Phi(x_j), \ i = 1, \ ..., \ m$$
 (6)

将公式(5)两边同乘 $\Phi(x_l)^{\mathrm{T}}$,得到

 $\lambda \Phi(x_l)^{\mathrm{T}} S_{\mathrm{T}}^{\Phi} W_i^{\Phi} = \Phi(x_l)^{\mathrm{T}} S_{\mathrm{B}}^{\Phi} W_i^{\Phi}, l = 1, 2, ..., N$ (7) 公式(5) 与公式(7) 有相同的特征向量解。

F 空间的维数非常高,加上 Φ是一个隐式的映射,使得 数值计算异常困难。为了便于数值计算,将上述公式表达为 F 空间中的内积形式。为此引入核函数^[4]

$$k(x, y) = (\Phi(x) \bullet \Phi(y)) \tag{8}$$

综合(4),(6),(8)整理得^[7]

$$\lambda K K \alpha = K M K \alpha \tag{9}$$

× N 的块对角阵, $M = (M_k)$, k = 1, ..., $c_{\circ} M_k$ 是元素均为 $1/n_k$ 的 $n_k \times n_k$ 的矩阵。

将 W_i^{Φ} 归一化得到: $\alpha_i^{\mathrm{T}} K \alpha_i = 1$ 。

对于测试数据 z_i 它在 F 空间中的像 $\Phi(z)$ 在 W_i^{Φ} 上的投影为

$$(W_i^{\Phi} \bullet \Phi(z)) = \sum_{j=1}^{N} \mathfrak{a}_{ij} (\Phi(x_j) \bullet \Phi(z)) =$$
$$\sum_{j=1}^{N} \mathfrak{a}_{ij} k(x_j, z)$$
(10)

总结上述讨论, GDA 的主要计算步骤如下。

(1) 在 F 空间要找到 Fisher 判决的最优投影方向等价于 求公式(5) 的特征向量解:

(2) 将求解公式(5)的 W° 转化为求解公式(9) 的 α ;

(3) 对于任意测试数据 z, 根据公式(10) 求其在 F 空间中
 ₩⁹ 方向上的投影。

可以看出,不必给出显示的非线性映射 Φ,所有的运算 都是通过输入空间中定义的内积核函数来完成的,这就是所 谓的核技巧。不同的核函数表示不同的非线性映射 Φ。常用 的核函数有高斯核,多项式核和 Sigmoid 核^{4]}。

3 广义判别分析与核主成份分析(KPCA)的 联系

文献[10]提出了 KFD(kernel fisher discriminant)的框架: 先基于 KPCA 的 m 个主成分得到空间 R^m , 然后在此空间进行线性判别分析。该文使用标准 Fisher 准则函数

 $J(W^{\Phi}) = \arg \max_{W^{\Phi}} (|(W^{\Phi})^{T} S_{B}^{\Phi} W^{\Phi}| / |(W^{\Phi})^{T} S_{T}^{\Phi} W^{\Phi}|)$ 在本小节中,我们使用扩展的准则函数(2),作为上述框架 的补充。

由公式(5)知:若 S^{P} 可逆,则 W_i^{O} 就是 $(S^{\text{P}})^{-1}S^{\text{O}}_{\text{B}}$ 关于特征值 λ_i 的特征向量。实际问题中 S^{P} 经常不可逆。为此有许多正则化方法来解决这个问题。

核 PCA 方法则是找到一组最优的投影方向最大化总散 度矩阵,即

$$J_{PCA}(W_{\circ\mu}^{\Phi}) = \arg \max_{W^{\Phi}} (|(W^{\Phi})^{T} S_{B}^{\Phi} W^{\Phi}| =$$
$$\arg \max_{W^{\Phi}} \frac{|(W^{\Phi})^{T} S_{T}^{\Phi} W^{\Phi}|}{|(W^{\Phi})^{T} W^{\Phi}|}$$

其中 | $(W^{\circ})^{T}W^{\circ}$ |= $I_{\circ}W^{\circ}$ 的列向量是互相正交的投影向 量,它是 S_{T}° 关于非零特征值 λ 的特征向量。由此看出,核 PCA 找到的投影方向也去除了 S_{T}° 的零空间,并且特征值 λ 的 大小代表了样本集方差的大小。设 $(\beta_{1}, \beta_{2}, ..., \beta_{m})$ 是 S_{T}° 相 应的 m 个最大特征值对应的特征向量, $m = \operatorname{rank}(S_{T}^{\circ})$ 。对于 输入空间的任意x, KPCA 变换后可以得到: $y = P^{T} \Phi(x)$,y= $(y_{1}, y_{2}, ..., y_{m})^{T}$, $P = (\beta_{1}, \beta_{2}, ..., \beta_{m})$ 。定义特征空间 F 中的子空间 $\Theta = \operatorname{span}{\beta_{1}, \beta_{2}, ..., \beta_{m}}$,它的正交补空间记 为 Θ^{\perp} 。实际上, Θ^{\perp} 就是 S_{T}° 的零空间。因此 $S_{T}^{\circ}W_{\Theta^{\perp}}^{\circ} = 0$, F= $\Theta + \Theta^{\perp}$ 。于是, F 空间中的扩展的 Fisher 准则函数为

$$J(W^{\Phi}) = \frac{(W^{\Phi})^{\mathrm{T}} S_{\mathrm{B}}^{\Phi} W^{\Phi}}{(W^{\Phi})^{\mathrm{T}} S_{\mathrm{B}}^{\Phi} W^{\Phi}}$$

其中, K, 为N × N 的矩阵, K_{ij} = ($\Phi(x_i)$ ・ $\Phi(x_i)$); M.为N $J(W = (W^{\phi})^T S^{\phi} W^{\phi} = (W^{\phi})^T S^{\phi} W^{\phi}$ (1994-2010 China Academic Journal Electronic Publishing House. All rights reserved. http://www.cnki.net

$$\frac{\left(W_{\Theta}^{\Phi}+W_{\Theta}^{\Phi\perp}\right)^{\mathrm{T}}S_{\mathrm{B}}^{\Phi}\left(W_{\Theta}^{\Phi}+W_{\Theta}^{\Phi\perp}\right)}{\left(W_{\Theta}^{\Phi}+W_{\Theta}^{\Phi\perp}\right)^{\mathrm{T}}S_{\mathrm{T}}^{\Phi}\left(W_{\Theta}^{\Phi}+W_{\Theta}^{\Phi\perp}\right)}$$
(11)

由干

$$(W_{\oplus}^{\oplus} + W_{\oplus}^{\oplus} \perp)^{\mathrm{T}} S_{\mathrm{T}}^{\oplus} (W_{\oplus}^{\oplus} + W_{\oplus}^{\oplus} \perp) = (W_{\oplus}^{\oplus})^{\mathrm{T}} S_{\mathrm{T}}^{\oplus} W_{\oplus}^{\oplus} + 2 (W_{\oplus}^{\oplus} \perp)^{\mathrm{T}} S_{\mathrm{T}}^{\oplus} W_{\oplus}^{\oplus} + (W_{\oplus}^{\oplus} \perp)^{\mathrm{T}} S_{\mathrm{T}}^{\oplus} W_{\oplus}^{\oplus} \perp =$$

 $(W_{\Theta}^{\Phi})^{\mathrm{T}} S_{\mathrm{T}}^{\Phi} W_{\Theta}^{\Phi}$

并且可以证明 S_T° 的零空间必然是 S_B° 的零空间 $^{(1)}$,从而(11)式的分子简化为: $(W_{\circ}^{\circ})^{T} S_B^{\circ} W_{\circ}^{\circ}$ 。所以公式(11)可以写成

$$J(W^{\Phi}) = (W^{\Phi}_{\Theta})^{\mathrm{T}} S^{\Phi}_{\mathrm{B}}(W^{\Phi}_{\Theta}) / (W^{\Phi}_{\Theta})^{\mathrm{T}} S^{\Phi}_{\mathrm{T}}(W^{\Phi}_{\Theta})$$
(12)

此时可以得出,基于扩展的 Fisher 准则函数的广义判别 分析依然可以纳入到文献[10]所提的框架中:在核 PCA 张 成的子空间中进行 LDA 可以得到核 Fisher 的投影向量,前 提是去掉了 S_1° 的零空间。

4 实验结果与分析

将 GDA 算法用于恒星、星系和类星体的自动分类。实验数据来自 SLOAN 的 DR2 的光谱数据库,其中恒星有1878条,星系2529条,类星体3026条。光谱波长范围:380~900 nm,步长取为1 nm,插值后每条光谱521个点。将整条光谱作为输入参与 GDA 的运算。

为了确定核函数的参数,恒星、星系和类星体光谱各取 1 000 条作为交叉验证样本集(训练样本集),剩下的作为测 试数据集。实验分别对比了线性判别分析(LDA)和GDA 算 法的结果。其中GDA 的核函数取为高斯核函数: $k(x, Y) = \exp(-||x - y||^2/\sigma^2)$ 。因为有三类数据,所以投影方向有 2 个。

图1所示的是三类数据基于 LDA 在第1、第2个投影方 向上的投影结果。图2所示的为三类数据基于 GDA 算法在 第1、第2个投影方向上的投影结果。可以看到,对于 LDA, 三类数据基本也能分开,但是数据的散度较大;对于 GDA, 三类数据的分布更加内敛,散度相对较小。可见 GDA 使得 不同类别样本的特征差别加大,而同类样本的特征差别减 小,因而 GDA 所提取的特征比 LDA 的特征更接易于分类。





用 LDA 和 GDA 分别提取特征后,用最近邻分类器 (KNN)完成分类。GDA 对于核宽有一定的敏感度,核宽不 同,分类的正确率相差较大。图 3 所示是在训练样本集上 GDA 在不同的核宽下的分类表现。实验按照 5 倍交叉验证 完成,核宽范围;0-01~10。可以看出,核宽较小时,分类正 确率高;核宽增大,分类正确率降低。最好的结果在 ∞= 0.2 时出现,为0.96。相比 LDA 的0.909 可以看出,选择合适的 核宽, GDA 的分类性能远高于 LDA。



Fig. 2 The projection of training data based on GDA



Fig. 3 Correct classification rate based with GDA on different kernel width





同时比较了分别以 PCA 和 KPCA 作为特征提取(KNN 为分类器)在交叉验证样本集(训练样本集)上的分类表现, 如图 4 所示。PCA 和 KPCA 分别取了 2~ 10 个主分量(其中 KPCA 的核宽为 2 0)。可以看出,当取 2 个主分量时,基于 PCA 的分类正确率为 0 773,基于 KPCA 的分类正确率为 0 778。这种分类结果比基于 LDA 和 GDA 的方法差很多, 但是随着维数的增大,基于 PCA 和 KPCA 方法的分类正确 率提升很快,PCA 取六维、KPCA 取五维时,分类正确率已 超过了 0 9。基于 KPCA 的分类效果比 PCA 要稍好一些。

对于含有4.433条光谱的测试数据集, LDA 的测试光谱

1962

分类的正确率为0 852, 其中恒星的分类正确率为0 952, 类 星体为0 815, 星系为0 85。GDA(σ= 0 2)的测试正确率为 0 936, 其中: 恒星的分类正确率为0 975, 类星体为0 94, 星系为0 906 5(见表1)。

Table 1 The correct classification rate based on LDA and GDA

输入 类别	LDA 分类正确率			GDA 分类正确率		
	Star	Quasar	Galaxy	Star	Quasar	Galaxy
Star	0 952	0 032	0 016	0 975	0. 025	0
Qu asa r	0 113	0 815	0 077	0 028	0.94	0 032
Galaxy	0 03	0 12	0 85	0 009 8	0 083 7	0 906 5

由表 1 看出, 经过线性 LDA 变换, Quasar 容易错分为 Star; 而 Quasar 和 Galaxy 也容易互相错分; Star 和 Galaxy 不容易互相错分。经过非线性的 GDA 变换后, Star, Quasar 和 Galaxy 的识别率都有增大, 以 Quasar 的识别率提高最多, 比 LDA 的识别率增加了近 14%, 说明 Quasar 对于基于核的 线性判别方法非常敏感; Galaxy 的识别率有所上升, 它还是 容易被错分为 Quasar; Star 的识别率上升最小; Quasar 和 Star 在非线性变换下不容易错分了; Star 和 Galaxy 的可分 性也进一步增强。这充分体现了基于核的判别分析的方法的 优越性。

5 总 结

近年来核方法在众多学科都得到广泛应用。本节主要讨 论了基于核的广义判别分析(GDA)算法在光谱分类中的应 用,并将基于扩展的Fisher 准则函数也纳入到 KFD(kernel fisher discriminant)算法的框架中,即:在核 PCA 张成的子 空间中进行线性判别分析(LDA)可以得到核 Fisher 的投影 向量,前提是去掉 Sf 的零空间。实验对比了 LDA,GDA, PCA,KPCA 算法对于恒星、星系和类星体的光谱分类表 现。结果是基于 GDA 的算法对于这三种类型光谱的分类正 确率最高,LDA 次之;尽管 KPCA 也是基于核的方法,但是 取低维时其效果不好,甚至低于 LDA,基于 PCA 的分类效 果最差(低维时)。这也说明有监督的方法在光谱分类问题上 比无监督的方法具有优势。

参考文献

- [1] QIN Dong mei, HU Zhar yi, ZHAO Yong heng(覃冬梅, 胡占义, 赵永恒). Spectroscopy and Spectral Analysis(光谱学与光谱分析), 2003, 23(1): 182.
- [2] Cabanac R A, de Lapparent V, Hickson P. Astronomy and Astrophysics, 2002, 389: 1090.
- [3] Weaver W B, et al. Astrophysical Journal, 1997, 487: 847.
- [4] Vladimir N Vapnik. Translated by ZHANG Xue gong(张学工译). Nature of Statistical Learning Theory (统计学习理论的本质). Bei jing: Tsinghua University Press(北京:清华大学出版社), 2000.
- [5] Scholk opf B, Smola A, Muller K R. Neural Computation, 1998, 10(5): 1299.
- [6] Mika S, Ratsch G, Weston J, et al. IEEE International Workshop Neural Networks for Signal Processing IX, 1999, (Aug.): 41.
- [7] Baudat G, Anouar F. Neural Computation, 2000, 12(10): 2385.
- [8] Richard O Du da, Peter E Hart, David G Stork. Transted by LI Hong dong, YAO Tiar xiang, et al(李宏东,姚天翔, 等译). Pattern Classification(Second Edition)(模式识别,第2版). Beijing: China Machine Press(北京: 机械工业出版社), 2003.
- [9] Chen Li fen, et al. Pattern Recognition, 2000, 33(10): 1713.
- [10] Yang Jian, Frangi A F, et al. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2005, 27(2): 230.
- [11] Huang R, Liu Q S, Lu H Q, et al. Solving the Small Sample Size Problem of LDA. Proceedings of 16th International Conference on Pattern Recognition, 2002, 3: 29C32.

Spectra Classification Based on Generalized Discriminant Analysis

XU Xin^{1,2}, YANG Jirr fu¹, WU Fu chao¹, ZHAO Yong heng²

- 1. National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing 100080, China
- 2. National Astronomical Observatories, Chinese Academy of Sciences, Beijing 100012, China

Abstract A kernel based generalized discriminant analysis (GDA) technique is proposed for the classification of stars, galaxies, and quasars. GDA combines the LDA algorithm with kernel trick, and samples are projected by nonlinear mapping onto the fear ture space F with high dimensions, and then LDA is conducted in F. Also, it could be inferred that GDA which combines the extension of Fisher's criterion with kernel trick is complementary to kernel Fisher discriminant framework. LDA, GDA, PCA and KPCA were experimentally compared with these three different kinds of spectra. Among these four techniques, GDA obtains the best result, followed by LDA, and PCA is the worst. Although KPCA is also a kernel based technique, its performance is not satisfactory if the selected number of the principal components is small, and in some cases, it appears even worse than LDA, a nor kernel based technique.

Keywords Spectra classification; Generalized discriminant analysis; Linear discriminant analysis; Kernel principal component analysis

(Received Aug. 8, 2005; accepted Nov. 8, 2005)