文章编号: 1672-8785(2013)06-0039-06

多指标混合逐次投影高光谱 图像降维算法

郑思远 李智勇 周石琳 王亮亮 (国防科技大学电子科学与工程学院,湖南长沙 410073)

摘 要: 降维对于高光谱图像解译具有重要意义。基于二阶统计量分析的经典主成分 分析方法在降维过程中会丢失小目标信号。为解决这一问题,本文中引入高阶统计量 作为投影指标对主成分分析方法进行拓展,提出了一种基于不同统计量描述的混合逐 次投影的高光谱图像降维算法。该方法在保持主成分分析优点的同时,有效结合了非 正交向量投影的特点,可以在降维后的低维空间中保留异常信号成分。真实高光谱图 像数据的实验结果证明,该方法相对于主成分分析可以提取更加完整的低维信号子空 间。

关键词: 高光谱图像; 降维; 主成分分析; 投影追踪; 高阶统计量

中图分类号: TP751.1 文献标识码: A DOI: 10.3969/j.issn.1672-8785.2013.06.08

Hyperspectral Image Dimension Reduction Algorithm Based on Multi–index Successive Projection

ZHENG Si-yuan, LI Zhi-yong, ZHOU Shi-lin, WANG Liang-liang (College of Electronic Science and Engineering, National University of Defense Technology, Changsha 410073, China)

Abstract: The dimension reduction is of significance to the interpretation of hyperspectral images. The classical Principal Component Analysis (PCA) method based on Second-order Statistics (SOS) may lose the signals of small targets during dimension reduction. To solve this problem, the PCA method is extended by using high-order statistics as projection indexes and a hyperspectral image dimension reduction method based on multi-index successive projection is proposed. The method incorporates the features of the non-orthogonal projection while retaining the advantages of the PCA method. It can keep the abnormal signal components in the dimension-reduced space. The experimental result of the real hyperspectral image data shows that this method can extract more complete signal subspace than the PCA method.

Key words: hyperspectral image; dimension reduction; principal component analysis; projection pursuit; high–order statistics

收稿日期: 2013-05-03

基金项目:国家自然科学基金项目(40901216)

作者简介:郑思远(1989-),男,新疆哈密人,硕士研究生,主要研究方向为模式识别和高光谱图像处理。 E-mail: nick.zsy@gmail.com

0 引言

高光谱遥感技术是当前遥感领域的前沿技 术。与传统遥感图像不同,高光谱图像由地物几 十个甚至上百个相邻光谱波段的图像构成, 包含 丰富的光谱信息。然而其高维数据空间内部大部 分是空的,数据常集中在低维结构中^[1]。而且高 光谱图像的高维数据特性也带来了数据量大、 数据冗余度高、计算复杂等问题。因此,通过合 适的降维方法将高维数据映射到低维空间对后 续的高光谱图像智能解译具有重要意义。主成分 分析 (principal components analysis, PCA)^[2] 是一种 基于数据二阶统计特性分析的线性变换降维算 法,在高光谱图像处理中得到广泛应用。然而, 在实际的高光谱图像处理中, PCA 方法对于一 些出现概率较低的异常像元、小目标不敏感, 在 降维过程中会丢失这部分信息^[3-5]。近些年基 于偏度 (skewness) 和峰度 (kurtosis) 等高阶统计量 的研究工作在高光谱图像处理中越来越受到重 视 [6-8]。信号高阶统计特性与高光谱数据特点 紧密结合,反映了该图像的数据分布与高斯分 布的偏离程度。其值越大,含有的异常信息就越 多,越有利于小目标、小类别信息的保留。

本文提出一种结合方差与高阶统计量作为 优化指标的混合逐次投影降维方法,以解决使 用 PCA 降维时容易丢失小目标信号的问题。第 二节从逐次投影的角度对 PCA 进行分析;第三 节介绍采用高阶统计量进行投影的方法,并从逐 次投影的角度对 PCA 方法进行扩展,提出多指 标混合逐次投影 (Hybrid Indexes for Iterative Projection,HIIP)算法,并在第四节中结合一种基于 多元数据样本偏度的正态性检验方法来自动确 定子空间维数。最后采用真实高光谱数据对本文 算法进行验证,并与 PCA 方法进行比较分析。

1 PCA

PCA 主要是通过对协方差矩阵进行特征分解,得出数据的主成分 (即特征向量) 与它们的权值 (即特征值)。设高光谱图像中 L 维像元光谱矢量观测样本集为 X=x_i,i=1,L,N, N 为像元

 ${\rm Infrared}~({\rm monthly})/{\rm Vol.34},~{\rm No.6},~{\rm Jun}~2013$

样本数量。利用 PCA 方法对高光谱图像进行降 维,首先对其进行中心化(去除均值)处理:

$$Y = X - \mu \cdot l^T \tag{1}$$

式中, $\mu = \frac{1}{N} \sum_{i=1}^{n} x_i$ 为样本均值, 1为分量均为 1 的 L 维列向量。由 PCA 方法进行降维的具体步 骤如下:

(1) 计算的协方差矩阵 $\sum_{N=1}^{\infty} = \frac{1}{N-1} Y Y^T$;

(2) 对 Σ 进行特征值分解, 将其全部特征值 按 $\lambda_1 \ge \lambda_2 \ge L \ge \lambda_L$ 的顺序排列, 各特征值对应 的特征向量为 e_i ;

(3) 选出前 k 个最大特征值所对应的特征向量 {*e_l*,L,*e_k*} 作为图像主成分,构成低维子空间的基向量组 *E_k*,对观测数据进行正交变换,得

$$Z = E_k^T Y \tag{2}$$

则可求得变换后样本数据集的协方差阵为

$$\hat{\sum}_{z} = \frac{1}{N-1} \sum_{i=1}^{n} Z Z^{T}$$
$$= E_{k}^{T} \hat{\sum} E_{k}^{T} = \begin{pmatrix} \lambda_{1} & L & 0\\ M & O & M\\ 0 & L & \lambda_{k} \end{pmatrix}$$
(3)

从结果可以看出,变换后低维空间中的数 据各分量是不相关的,而且各波段分量的方差 是其对应的特征值λ_i。通过去掉最小特征值所 对应的成分,可达到降维的目的。从投影追踪^[9] 的观点来看,变换矩阵中*E_k*的特征向量都是根 据方差最大化的欧氏空间投影方向确定的。主成 分分析的过程实际上也可以被看成是以二阶统 计量方差为投影指标的逐次投影寻踪过程,其 第一主成分的估计为

$$w_{1} = \arg \max_{||w||=1} E\{(w^{T}Y)^{2}\}$$

= $\arg \max_{||w||=1} E\{w^{T}YY^{T}w\}$ (4)

由瑞利商定理可知, w₁ 即为 *E*[*YY^T*] 的最 大特征值所对应的特征向量。这与之前 PCA 的 结果是一致的,并且其各主成分投影方向是相 互正交的。假设已得到前 k-1 个主成分 *W*_{k-1} = [*w*₁, *Lw*_{*k*-1}],为了得到第 k 个主成分,必须先从 Y 中减去前面的 k-1 个主成分:

$$Y_{k-1} = Y - \sum_{i=1}^{k-1} w_i w_i^T Y = Y - U_{k-1} Y \qquad (5)$$

式中, $U_{k-1}=W_{k-1}W_{k-1}^{T}$ 为投影矩阵。把求得的 第 k 个主成分带入数据集,得到新的数据集,继 续寻找主成分。

$$w_k = \arg \max_{||w||=1} var\{w^T \hat{Y}_{k-1}\}$$
(6)

2 混合逐次投影降维

2.1 基于高阶统计量的最佳投影方向

从逐次投影过程理解 PCA, 其核心就是以 方差为指标的相互正交投影矢量寻找过程。实 际上利用数据更高阶的统计信息进行投影可得 到具有不同性质的投影方向。

偏度与峰度由随机变量的三阶与四阶中心 矩定义,在统计学中分别被用来衡量概率分 布相对于正态分布的偏斜度以及平坦度。对 于中心化数据 $Y = \{y_i\}_{i=1}^N$, 令 $z = w^T Y = (w^T y_1, L, w^T y_N)^T = (z_1, L, z_N)^T$,其三阶、四阶 统计量偏度与峰度可表示为

$$k_{3}(z) = E[(w^{T}Y)^{3}] = \frac{1}{N} \sum_{i=1}^{N} z_{i}^{3}$$
(7)
$$k_{4}(z) = E[(w^{T}Y)^{4}] = \frac{1}{N} \sum_{i=1}^{N} z_{i}^{4}$$

可类似定义更高阶统计量:

$$k_k(z) = E[(w^T Y)^k] = \frac{1}{N} \sum_{i=1}^N z_i^k, k \ge 3 \qquad (8)$$

相对于二阶统计量方差来讲,高阶统计量对异 常数据更为敏感。在高光谱图像处理过程中, 高阶统计量可以用来寻找大面积背景下的小物 体。与使用方差 *E*[*w^TY*)²] 作为投影指标的 PCA 不同,当将投影指标换作高阶统计量 *k*_k(*z*) 时, 该最优化问题就可表示为

$$w^* = \max_{||w||=1} \{\frac{1}{N} \sum_{i=1}^{N} z_i^k\}$$

http://journal.sitp.ac.cn/hw

$$= \max_{||w||=1} \{ \frac{1}{N} \sum_{i=1}^{N} w^{T} y_{i} (y_{i}^{T} w)^{k-2} y_{i}^{T} w \}$$
(9)

Ren 采用拉格朗日乘数法对上述问题进行了分析^[10],并构造了目标函数:

$$J(w) = E[w^T Y(Y^T w)^{k-2} Y^T w] - \lambda(w^T w - 1)$$
(10)

令 $\frac{\partial J(w)}{\partial w} = 0$, 可得到最优解表示:

$$(E[Y(Y^T w)^{k-2} Y^T] - \lambda' I)w = 0$$
(11)

E[*Y*(*Y^Tw*)^{*k*-2}*Y^T*]的最大特征值所对应的特征向 量即为所求的最佳投影方向 *w*^{*}。这里给出一种 求解 *w*^{*} 的迭代算法:

(1) 初始化随机投影矢量 $w^{(0)}$, n=0;

(2) 计算 $E[Y(Y^Tw)^{k-2}Y^T]$ 及其最大特征值 所对应的特征向量 $v^{(n)}$;

(3) 令 $w^{(n+1)} = v^{(n)}$, 计算 $w^{(n+1)} = w^{(n)}$ 之 间的夹角 θ 。如果 $\theta < \varepsilon$,即 $w^{(n+1)} = w^{(n)}$ 趋于 一致,那么 $w^* = v^{(n)}$;否则, $n \leftarrow n+1$,转b 继续执行。

2.2 HIIP

不同的投影指标,其投影方向有不同的统 计含义。如果使用方差作为投影指标,则得到的 投影方向上的数据分布信息量最大,像元能量 分布最集中。如果以偏度作为衡量准则,投影方 向上的数据直方图分布的偏斜性最大,这在一 定程度上说明该方向上可能含有的异常信息较 多。而如果选择峰度,则会导致投影方向上的数 据分布具有长拖尾性质。

使用 PCA 对高光谱图像进行降维时,只需 对其二阶统计量信息进行优化,可以较好地保 存图像中的大部分背景信息,但是容易丢失一 些出现概率较小的异常或者小目标信号。文献 ^[9] 提出了以某一种高阶统计量作为优化指标的 迭代投影矢量生成算法。该方法可以较好地提 取小目标信号,但是会丢失一部分区域面积中 的较大目标的信息。在很多情况下,我们希望在 对图像进行降维的同时能保存好原始数据中的 背景以及异常信息。为了达到上述目的,一种可 行的方案就是结合不同优化指标下的投影方向 生成投影空间,以便充分利用图像中的信息。我

INFRARED (MONTHLY)/VOL.34, NO.6, JUN 2013

们提出一种不同优化指标混合的逐次投影算法 HIIP 生成信号子空间,其具体过程如下:

(1) 对高光谱数据样本集 X 进行去均值处理,得到中心化数据 Y。

(2) 采用不同阶数的统计量 $\kappa_k(z)$ 作为投影 指标,在Y上求解其最佳投影方向;本文选取 k=2、3、4,即分别为方差、偏度以及峰度, 所得到的最佳投影方向分别为 w_1 、 w_2 、 w_3 。 记为 $W_{(1)} = [w_1, w_2, w_3]$ 。

(3) 令 $P_{W_{(1)}} = W_{(1)}(W_{(1)}^T W_{(1)})^{-1} W_{(1)}^T$, $P_{W_{(1)}}^{\perp} = I - P_{W_{(1)}}$,可得到去除 $W_{(1)}$ 投影成分的新数据 集 Y^1

$$Y^{1} = P_{W(1)}^{\perp} Y$$
 (12)

(4) 用 Y^1 替换 Y,执行步骤 2,产生新的 投影方向组 $W_{(2)}$,则 $Y^2 = P_{W_{(2)}}^{\perp}Y$;假设已经 得到前 k-1 组投影方向组 $W_{(1)}$, $LW_{(k-1)}$,则 数据第 k 级成分 $Y^k = P_{W_{(k-1)}}^{\perp}Y^{k-1}$,之后在 Y^k 上执行 2,得到 $W_{(k)}$ 。

(5) 在达到可满足要求的低维子空间维数 后停止迭代,记录子空间的基向量组为 $W = [W_{(1)}, LW_{(k)}]$,所得低维信号空间即为 $\hat{S}_k = range\{W_{(1)}, LW_{(k)}\}$ 降维数据,为

$$Z = W^T Y \tag{13}$$

另外需要注意的是, 在步骤 4 中生成 Y^k 还 可以通过用原数据 Y 直接投影获得, 即 $Y^k = P_{[W_{(1),LW_{(k-1)}]}}^{\perp} Y$ 。上述过程可以看作是对主成分 分析的一种推广, 实际上就是在 PCA 所寻找的 最大方差投影的基础上增加高阶投影方向的信 息。从几何角度的观点来看, PCA 生成的各级 主成分投影方向是相互正交的。这种相互正交 的投影方式可以在最大程度上保存数据原有的 主要信息。但在高光谱图像中存在的一些异常 信号, 完全正交的投影并不是最佳投影方式。而 基于高阶统计量的投影方向并不一定与主成分 方向正交, 这种非正交方向在一定程度上弥补 了各级正交主成分的不足。

3 子空间维数确定

对于高光谱图像降维,需要考虑的另一个问题就是如何确定最终子空间的维数。在采用 逐次投影的方法时,子空间维数的确定实际上 就是迭代终止条件的选择。考虑到图像由信号 S 与噪声 N 构成,即

$$Y = S + N \tag{14}$$

如果认为提取的子空间完全包含了原始图像的 全部信息,则可认为投影到其补空间的数据为 不含任何信息的正态性噪声,即

$$Y^{k} = P_{\hat{S}_{\mu}}^{\perp} Y \approx P_{\hat{S}_{\mu}}^{\perp} N \tag{15}$$

因此,我们考虑以下的二元假设检验:

 $H_0: \{Y^k \text{ is normal}\}$

 $H_1:\{Y^k \text{ is non-normal}\}$ (16)

当假设 H₀ 成立时,即可认为降维子空间已 满足要求,停止迭代过程。



图 1 多指标逐次混合投影算法的流程图

对于 (16) 中的正态性检验问题, Móri 定义 了一种用来检验多元数据正态性分布的多元随 机变量的偏度形式 ^[11]:

$$\beta_{1,d}^{\%}(y) = E(||x||^2 x) \tag{17}$$

式中, $x = \sum^{-1/2} (z - \mu)$ 。对 Y^k 进行白化之后, $\beta_{1d}^{\%}$ 的数据样本形式为

$$b_{1,d}^{\%} = \frac{1}{N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} ||y_i||^2 y_i^T y_j||y_j||^2$$
(18)

当 N → ∞ 时, $b_{1,d}^{\%}$ 的极限分布趋近于正态分布 [12]:

$$Nb_{1,d}^{\%} \longrightarrow 2(d+2)\chi_d^2$$
 (19)

在给定显著性水平 α 下, 当 $Nb_{1,d}^{\%}/2(d+2)$ 的值 超过 χ_{a}^{2} 的 1 – α 分位数 $\chi_{d,1-\alpha}^{2}$ 时,则拒绝式 (16) 中的原假设 H_{0} ;当接受 H_{0} 时,整个迭代过程 结束,即可得到逐次混合投影的子空间。结合前 面投影过程的介绍,逐次混合投影算法的流程 图如图 1。

4 混合逐次投影降维



图 2 PHI 数据的第7波段的场景图

我们采用由上海技物所研制的推扫式成像 光谱仪生成的高光谱图像对算法进行验证。PHI 数据由机载平台在实际飞行中获得,剔除部分信 噪比较低、图像模糊的波段。实际参与处理的波 段有 123 个,图像大小为 220×90。图 2 为主要场 景示意图,主要地物包括机场跑道、沙地以及 7 台测试用车辆。

本文采用 Matlab 作为实验平台实现 HIIP 算法,并将实验结果与 PCA 方法进行对比。图 3 为 PCA 以及 HIIP 的投影成分分量图。实验中,

HIIP 算法正态性检验的显著性水平参数设置为 $\alpha = 10^{-4}$,用 HIIP 得到 3 个基向量组,共9 个分 量。图 3 (a)~(e)分别是 HIIP 得到的第 1、3、 5、7、9 成分的影像,(f)~(j)为 PCA 得到的与 之相对应的第 1、3、5、7、9 成分影像。从 投影成分的结果中可以看出,PCA 中前几个主 成分包含了图像的大部分背景信息,其后面几 个成分中几乎没有有效信息。这说明在通过主 成分变换后,大部分数据都集中在前几个主成分 方向上,而其他与主成分正交的方向上的信息 量趋于 0。但是从 (c)~(e)的结果来看,以偏度与 峰度作为投影指标时,在非正交投影方向上仍 含有一些重要的小目标信息。



图 3 HIIP 与 PCA 投影成分的影像

为进一步比较 PCA 与 HIIP,我们将原始 图像数据分别投影到 HIIP 得到的信号子空间以 及由 PCA 前 9 个主成分构成的子空间的补空间 中:

$$\bar{Z} = P_{\hat{S}}^{\perp} Y \tag{20}$$

投影到补空间的信号即可认为是降维过程中损失的信号成分。则各像素点损失能量可定义为

$$E_{loss} = ||\bar{z}_i||_2^2 = ||P_{\hat{S}}^{\perp} y_i||^2 \tag{21}$$

图 4 中分别展示了两种算法损失能量的分布图, 其中图 4(a) 与 4(b) 分别为 PCA 与 HIIP 的能量 损失。从中可明显看出 PCA 方法中仍有关于车辆目标的信号成分损失。而与 PCA 相比,采用 HIIP 时,小目标的能量损失明显减少了。



(b)

图 4 损失的能量分布示意图。(a) PCA, (b) HIIP



图 5 补空间投影信号能量正态的拟合示意 图。(a)PCA, (b)HIIP

理想情况下,通过降维手段得到的子空间应 包含原始图像中的全部信号成分,而在其补空 间中为不含有任何信息的高斯噪声残余信号。 假设补空间投影信号为零均值高斯分布,则其 分量平方和形式服从分布。当样本数量足够多 时,其均方根近似服从正态分布。因此,通过衡量补空间投影信号能量均方根的正态性,可以从侧面反映低维信号投影空间的准确程度。图 5中(a)为PCA方法的残余信号能量正态拟合结果,(b)为HIIP残余能量的拟合结果。从图中可以看出,由PCA得到的子空间在其补空间中仍有非正态信号成分,使得残余能量形式明显偏离正态分布。逐次混合投影算法通过引入高阶统计量信息对主成分分析进行有效补充,得到了更完整的信号子空间。

5 结论

在高光谱图像的降维方法中,基于二阶统 计量分析的主成分分析方法可以有效保存图像 中的背景信息,但会丢失出现概率较低的异常 信号成分。本文从逐次投影追踪的角度对 PCA 方法进行了分析,并在此基础上通过引入基于 高阶统计量的投影方法对 PCA 进行了拓展。然 后根据此思想设计了一种 HIIP 的高光谱数据降 维方法。从实验结果来看,本文算法减少了 PCA 降维过程中的异常信号损失,与 PCA 相比,所 得到的信号子空间有效增加了对于小目标信号 的表示。今后还需要进一步研究如何将该方法 与高光谱图像异常检测、小目标检测以及分类 等应用相结合,以提高算法的性能。

参考文献

- Landgrebe D. Hyperspectral Image Data Analysis
 [J]. IEEE Signal Processing Magazine, 2002, 19(1): 17–28.
- [2] Bajorski P. Statistical Inference in PCA for Hyperspectral Images [J].Selected Topics in Signal Processing, 2011, 5(3): 438–445.
- [3] Kuybeda O, Malah D, Barzohar M. Rank estimation and Redundancy Reduction of High-dimensional Noisy Signals with Preservation of rare Vectors [J].IEEE Transactions on Signal Processing, 2007, 55(12): 5579–5592.
- [4] Prasad S, Bruce L M. Limitations of Principal Components Analysis for Hyperspectral Target Recognition [J]. *IEEE Geoscience and Remote Sensing Letters*, 2008,5(4): 625–629.

(下转第48页)