

Dimensionality Reduction for
Hyperspectral Remote Sensing Data

高光谱遥感数据降维

王雪松 程玉虎 孔毅 高阳 著



科学出版社

高光谱遥感数据降维

王雪松 程玉虎 孔毅 高阳 著

科学出版社

北京

内 容 简 介

高光谱数据降维是遥感数据在土地资源分析及应用的第一步,是人们获取遥感信息的一种重要手段。针对高光谱数据具有的高维数、非线性、数据量大、标记样本少等特性,利用机器学习、模式识别和遥感科学等多学科交叉的理论和方法,研究高光谱数据降维问题。

本书以稀疏表示、张量学习、迁移学习和深度学习为基础,系统阐述如何更好地进行非线性特征提取、如何针对高光谱数据的三维特性设计降维算法以及如何充分利用有限标记样本和大量未标记样本来提高降维效果等问题。各章节均涉及相关领域基础知识的介绍,能够为不同层次的读者与研究人员提供入门知识与参考信息。

本书可作为模式识别、数据挖掘、地理信息系统相关专业研究生的辅助教材,或作为机器学习工程师的参考书目。

图书在版编目(CIP)数据

高光谱遥感数据降维/王雪等著. —北京:科学出版社,2017.6

ISBN 978-7-03-053126-1

I. ①高… II. ①王… III. ①遥感图象-图象处理 IV. ①TP751

中国版本图书馆 CIP 数据核字(2017) 第 126782 号

责任编辑:惠 雪 曾佳佳/责任校对:李 影

责任印制:张 倩/封面设计:许 瑞

科学出版社 出版

北京东黄城根北街 16 号

邮政编码:100717

<http://www.sciencep.com>

文林印务有限公司 印刷

科学出版社发行 各地新华书店经销

*

2017 年 6 月第 一 版 开本:720 × 1000 1/16

2017 年 6 月第一次印刷 印张:11 1/2

字数:232 000

定价:69.00 元

(如有印装质量问题,我社负责调换)

前 言

随着航天技术、计算机技术、通信技术、信息处理技术的进步,现代空间遥感技术得到了空前发展,日益朝着“三全”(全天候、全天时和全球观测)、“三高”(高空间分辨率、高光谱频率和高时相分辨率)和“三多”(多传感器、多平台和多角度)方向迅猛发展。

与多光谱数据相比,利用高光谱数据(又称高光谱遥感数据、高光谱图像、高光谱遥感图像)进行地物识别和分类的优势主要体现在以下几个方面:①识别更多的地物,并且可以区分同类地物间的细微差异;②可选择的样本数据和光谱特征具有多样性和灵活性;③图谱合一。这些特点和优势,使得在多光谱数据中不可识别的地物,在高光谱数据中能够被很好地识别出来。因此,高光谱数据技术的应用具有广泛的背景和良好的前景,并与合成孔径雷达、激光探测与测距一起被视为今后最具发展前景的三种遥感信息获取技术。

理论上,随着高光谱数据光谱分辨率的提高,特征空间的维数也越来越高,因而表现不同地物类别的能力也随之不断提高。然而,实际应用中,现有的高光谱数据地物识别与分类算法却远远不能满足高光谱数据传感器技术的快速发展需求。高光谱数据分类主要面临下述三个问题:①“维数灾难”问题;②分类代价昂贵问题;③非线性可分问题。因此,如何既充分利用高维特征空间所提供的充足信息量,又解决特征空间维数过高的问题,是进行高光谱数据分析研究的关键。通常情况下,可以利用特征选择或特征提取等降维算法把数据从高维空间投影到低维空间中。

在国家自然科学基金项目(61273143、61472424)的资助下,本书从高光谱数据的特点出发,运用稀疏表示、张量学习、迁移学习以及深度学习等技术,研究高光谱数据的降维算法,共5部分11章。第1部分内容为第1~2章,主要为高光谱降维概述,综述内容包括:高光谱数据和降维算法的研究现状以及高光谱数据降维的相关研究基础。第2部分内容为第3~8章,旨在挖掘高光谱数据光谱信息的非负性、稀疏性和流形特性,包括:基于样本依赖排斥图的非负稀疏嵌入投影的高光谱数据降维、基于加权近邻保持嵌入的高光谱数据降维、遥感影像的半监督判别局部排列降维、基于块非负稀疏重构嵌入的高光谱数据降维、基于非负稀疏图的高光谱数据降维、基于非负稀疏半监督的高光谱数据降维。第3部分内容为第9章,围绕

高光谱数据的张量型降维展开,涉及张量表示、张量距离、高质量近邻图构建和张量型补丁校准框架等方法。第4部分内容为第10章,针对源高光谱数据和目标高光谱数据来自不同分布的问题展开,研究基于成对约束判别分析-非负稀疏散度的高光谱数据降维。第5部分内容为第11章,旨在利用深度学习挖掘高光谱数据的深层特征,研究基于样本依赖排斥图正则化自动编码器的高光谱数据降维。

由于时间仓促且作者水平有限,书中不当之处在所难免,恳切希望得到广大读者的批评和指正。

著 者

2017年2月于中国矿业大学

目 录

前言

第 1 章 高光谱研究概述	1
1.1 高光谱数据研究现状	3
1.1.1 高光谱数据发展现状	3
1.1.2 高光谱数据应用领域	4
1.1.3 高光谱数据降维的研究现状及存在问题	5
1.2 降维算法的研究现状	7
1.2.1 基于稀疏表示的降维算法	7
1.2.2 张量型降维算法	8
1.2.3 基于特征的迁移学习方法	8
1.3 本书主要研究方法	9
参考文献	11
第 2 章 高光谱数据降维研究基础	19
2.1 高光谱数据分析	19
2.1.1 高维数据的几何特征	19
2.1.2 高维数据的统计分布特征	22
2.2 高光谱数据特点	24
2.2.1 高光谱数据的数据模型	24
2.2.2 高光谱数据的空间相关性	25
2.2.3 高光谱数据的谱间相关性	27
2.2.4 Hughes 现象	28
2.3 高光谱数据的降维算法	29
2.3.1 特征选择方法	30
2.3.2 特征提取方法	30
2.4 高光谱数据降维的分类评价指标	45
2.4.1 混淆矩阵	46
2.4.2 整体分类精度	46
2.4.3 使用者、生产者及平均精度	46
2.4.4 Kappa 系数	47
2.5 高光谱数据降维的实验数据集	47

2.5.1	AVIRIS 高光谱数据	47
2.5.2	Hyperion 高光谱数据	50
2.5.3	ROSIS University 高光谱数据	51
2.5.4	ProSpecTIR 高光谱数据	52
2.6	本章小结	53
	参考文献	53
第 3 章	基于样本依赖排斥图的非负稀疏嵌入投影高光谱数据降维	57
3.1	基于样本依赖排斥图的非负稀疏嵌入投影	58
3.1.1	非负稀疏表示	59
3.1.2	样本依赖排斥图构建	59
3.1.3	低维嵌入投影	61
3.1.4	算法步骤	63
3.2	实验与分析	63
3.3	本章小结	67
	参考文献	68
第 4 章	基于加权近邻保持嵌入的高光谱数据降维	70
4.1	分布形变和加权距离	71
4.2	加权近邻保持嵌入	73
4.3	算法步骤	75
4.4	实验与分析	75
4.4.1	人工数据集	75
4.4.2	AVIRIS 高光谱遥感实验数据	77
4.5	本章小结	80
	参考文献	80
第 5 章	遥感影像的半监督判别局部排列降维	82
5.1	判别局部排列	83
5.2	基于图的半监督判别局部排列	84
5.3	实验与分析	86
5.4	本章小结	89
	参考文献	89
第 6 章	基于块非负稀疏重构嵌入的高光谱数据降维	91
6.1	块非负稀疏重构嵌入	92
6.1.1	非负稀疏表示	92
6.1.2	块非负稀疏表示	93
6.1.3	低维嵌入	93

6.1.4 算法步骤	95
6.2 实验与分析	95
6.2.1 高光谱实验数据	96
6.2.2 BNSRE 中字典块个数及稀疏权重矩阵分析	97
6.2.3 降维性能分析	99
6.2.4 讨论	100
6.3 本章小结	102
参考文献	102
第 7 章 基于非负稀疏图的高光谱数据降维	104
7.1 基于非负稀疏图的降维	105
7.1.1 问题描述	105
7.1.2 块非负稀疏表示	106
7.1.3 非负稀疏图构建	107
7.1.4 目标函数	109
7.2 算法步骤	110
7.3 实验与分析	111
7.4 本章小结	115
参考文献	115
第 8 章 基于非负稀疏半监督的高光谱数据降维	118
8.1 非负稀疏半监督降维算法	119
8.1.1 判别项	120
8.1.2 正则项	121
8.1.3 非负稀疏半监督最大间隔准则	122
8.2 实验与分析	123
8.3 本章小结	127
参考文献	128
第 9 章 基于高质张量近邻图和补丁校准的高光谱数据降维	130
9.1 基于高质张量近邻图和补丁校准的降维	131
9.1.1 高光谱数据光谱-空间信息的张量表示	132
9.1.2 张量距离	133
9.1.3 高质量近邻图	134
9.1.4 张量型补丁校准	135
9.2 实验与分析	138
9.2.1 参数分析	139
9.2.2 降维性能分析	142

9.3 本章小结	147
参考文献	148
第 10 章 基于成对约束判别分析-非负稀疏散度的高光谱数据降维	150
10.1 基于成对约束判别分析-非负稀疏散度的降维	151
10.1.1 问题描述	151
10.1.2 成对约束判别分析	153
10.1.3 非负稀疏散度准则	155
10.1.4 算法步骤	157
10.2 实验与分析	158
10.2.1 参数分析	158
10.2.2 对比实验	160
10.3 本章小结	163
参考文献	163
第 11 章 基于样本依赖排斥图正则化自动编码器的高光谱图像降维	167
11.1 基于样本依赖排斥图正则化自动编码器的降维	168
11.1.1 问题描述	168
11.1.2 样本依赖排斥图构建	169
11.1.3 基于样本依赖排斥图正则化自动编码器	169
11.2 算法步骤	172
11.3 实验与分析	172
11.4 本章小结	175
参考文献	175

第 1 章 高光谱研究概述

20 世纪 60 年代以来,人类在航空航天信息获取和卫星对地观测方面成绩斐然,卫星遥感技术快速发展。人类已迈向构建天地一体化对地观测系统的阶段,一种跨国家、跨组织的综合、持续、协同的分布式对地观测系统正在悄然形成,其目的是协调目前全球各自独立运行的各种监测平台、资源和网络,弥合系统之间的鸿沟,支持系统间的协同工作,逐步建设一个由多系统组成的综合、持续、协同的分布式对地观测系统,以保障监测和跟踪全球各个角落地球环境的变化,为全球性、国家性、地区性、部门性的环境、健康、灾害等社会公益事业的政策制定、决策与服务提供更快、更多、更好的数据与信息服务^[1]。

随着航天技术、计算机技术、通信技术、信息处理技术的进步,现代空间遥感技术得到了空前发展。高光谱数据技术是 20 世纪 80 年代以来综合对地观测的重要组成部分,也是国际对地观测技术竞争的关键点之一,朝着“三全”(全天候、全天时和全球观测)、“三高”(高空间分辨率、高光谱频率和高时相分辨率)和“三多”(多传感器、多平台和多角度)方向迅猛发展^[1,2]。

高光谱数据主要针对地物中的每一像元,利用紫外、可见光、近红外和中红外等大量电磁波波段,获取近似连续的地物光谱曲线(数据立方体)。高光谱数据具有“图谱合一”的特性,将反映地物反射特性的光谱波段信息和反映地物空间位置关系的图像信息结合在一起,可以更有效地对地物进行识别和分类,进而更有效地认识地物的有用特征信息、分布与变化规律。与多光谱数据相比,利用高光谱数据进行地物识别和分类的优势主要体现在以下几个方面^[3-5]。

1) 识别更多的地物,并且可以区分同类地物间的细微差异

高光谱数据获得光谱波段数多,一般是几十个或几百个,有的甚至达到几千个,这些地物光谱曲线在一定范围内是连续的、细微的,它能够真实地反映地物的本质特征信息、分布与变化的规律。传统的多光谱数据获取的光谱波段数只有几个,光谱分辨率低,导致其获取的光谱中存在“异物同谱”和“同谱异物”现象。由此可见,高光谱数据能够区分不同或同类地物间千差万别的光谱特征,大大提高地物的识别率。

2) 地物识别和分类时可选择的样本数据和光谱特征具有多样性和灵活性

高光谱数据波段数多且随着波段数的增加,数据量呈指数增加,包含丰富的信息。通过选择不同的样本和波段组合,可以获得不同的特征信息,为地物识别和分

类提供广泛的信息分析空间。

3) 图谱合一, 可以定量分析遥感信息

因为多光谱数据的光谱分辨率低, 所以其主要以定性分析为主, 部分定量分析的效果不好。而高光谱数据“图谱合一”的特性, 使获得的光谱曲线几乎和地面实测的同类地物的光谱曲线相同, 能够较好地提取地物的反射特性参量, 从而可以定量分析与提取地物的可分信息。

高光谱数据的这些特点和优势, 使得在多光谱数据中不可识别的地物, 在高光谱数据中能够被很好地识别出来。因此, 高光谱数据技术在环境监测、植被的精细分类、农作物的长势监测、城市热岛效应^[6]、地质岩矿的识别、地震灾害快速监测^[1]、油气微渗漏信息的监测、海洋水色定量监测^[2]等方面有良好的应用前景。高光谱数据技术与合成孔径雷达 (synthetic aperture radar, SAR)、激光探测与测距 (light detection and ranging, LiDAR) 一起被视为今后最具发展前景的三种遥感信息获取技术^[3]。

理论上, 随着高光谱数据光谱分辨率的提高, 特征空间的维数越来越高, 因而表现不同地物类别的能力也随之不断提高。然而, 实际应用中, 现有的高光谱数据地物识别与分类算法却远不能满足高光谱数据传感器技术的快速发展需求。这主要存在下面几个问题。

1) “维数灾难^[7]”问题

高光谱数据往往具有上百个波段, 数据量非常大。在提供丰富、详细信息的同时, 不同波段特别是相邻波段之间通常具有较强的相关性, 导致信息冗余量大。当所有波段全部用于分类时, 不仅效率低、计算量大, 而且波段数并不与精度成正比, 甚至波段数超过一定数目后整体分类精度反而会降低, 形成所谓 Hughes 现象^[7]。

2) 分类代价昂贵

在二维特征空间中, 500 个训练样本已足够多, 但在 200 维空间中则显得微不足道。尤其是针对基于统计学的分类器而言, 随着空间维数的增加, 待估计的统计参数也急剧提高, 若要获得好的整体分类精度, 就需要大量的训练样本, 这样分类时间也会随之增加。而高光谱数据分类的样本不易获得, 需要很多专家耗费较多的时间与精力, 代价非常高。

3) 非线性可分问题

由于多种原因, 如不同时间的大气条件、不同的采集系统状态、不同层次的土壤水分^[8]、不同的反射率和照明条件等环境下获取的同类地物的光谱曲线会有所不同, 所以高光谱数据在高维空间中具有非线性可分特点。

因此, 受上述三个问题的限制, 特征空间的维数不能太高。多光谱数据的特征空间维数虽然比较低, 但由于光谱分辨率较低, 不足以精确地表现不同的地物。而高光谱数据虽然光谱分辨率较高, 但由于波段数太多, 相应的特征空间维数太高,

带来了“维数灾难”、分类代价昂贵和非线性可分等问题。因此,如何既充分利用高维特征空间所提供的充足信息量,又解决特征空间维数过高的问题,是进行高光谱数据分析研究的关键。

通常情况下,利用特征选择和特征提取等降维算法把数据从高维空间投影到低维空间中。本书从高光谱数据的特点出发,运用非负稀疏表示、张量学习、知识迁移和深度学习理论,研究高光谱数据的降维算法,希望提高高光谱数据的整体分类精度、学习效率和泛化能力并减轻用户的负担。研究成果不但可以为解决高光谱数据降维问题提供新的分析设计方法和技术储备,而且可以进一步深化和丰富现有的机器学习、模式识别、遥感科学等理论。本书是遥感科学、机器学习、模式识别、计算机科学等多门学科有机交叉、新颖且富有挑战性的研究方向,具有重大理论意义和实际应用价值。

1.1 高光谱数据研究现状

1.1.1 高光谱数据发展现状

高光谱数据技术的发展主要分为两个大的方面。

1. 高光谱数据探测技术的发展

近 30 多年来,随着现代科技的发展,出现了航空遥感技术、卫星遥感技术等多种多样的探测技术,以及诸如微波、红外、雷达、激光、地磁等相关的各种遥感传感器。1983 年,第一幅高光谱数据由美国宇航喷气推进实验室 (Jet Propulsion Laboratory, JPL) 研发的世界上首台航空成像光谱仪 (aero imaging spectrometer-1, AIS-1) 拍摄^[9,10],标志着第一代成像光谱仪面世,并成功在矿物勘测、植被遥感、化学分析等方面应用。第一代成像光谱仪开创了“图谱合一”的高光谱数据技术的新时代。在此基础上,1987 年,第二代成像光谱仪面世,主要代表有 JPL 研制的航空可见光/红外光成像光谱仪 (airborne visible/infrared imaging spectrometer, AVIRIS)^[11,12]、美国研制的高光谱数字图像实验仪 (hyperspectral digital imagery collection experiment, HYDICE)^[13,14]、德国研制成功的反射式成像光谱仪 ROSIS、美国 EO-1 系统搭载的高光谱成像仪 Hyperion 等。这些成像光谱仪被成功应用到农业、生态、林业、天气、海洋、资源勘查等诸多领域。第三代成像光谱仪是傅里叶变换高光谱成像仪 (Fourier transform hyperspectral image, FTHSI)^[15,16]。

近几年,国内成像光谱仪的研制获得较大进展。在“八五”期间,我国研制了具有 64 个波段的可见光/近红外模块化机载成像光谱仪 (MAIS)。在国家 863 计划的支持下,“九五”期间,我国研制了具有 128 个波段模块化航空成像光谱仪 (OMIS) 和 224 个波段的机载推扫式高光谱成像仪 (PHI)^[16,17]。

2. 高光谱数据处理、分析的理论和方法

目前,高光谱数据分析方法主要有两个方向。第一个方向是基于光谱空间的分析方法。其基本原理是化学分析领域常用的光谱分析技术,主要是通过分析不同地物的光谱曲线表现出的不同光谱特征,来达到地物识别与分类的目的。这方面比较成熟的方法有光谱角填图技术^[18]、线性光谱分解技术^[19]、光谱匹配滤波技术^[20]、光谱特征匹配技术^[21]等。第二个方向是基于特征空间的分析方法。其基本思想是把组成光谱曲线的各光谱波段组成高维空间中的一个向量,进而用空间统计分析的方法分析不同地物在特征空间中的分布规律。两类方法各有优缺点:基于光谱空间的分析方法更直接,且不需要太多的地面先验知识,但由于各种因素的影响,其光谱曲线中往往会有许多噪声,给光谱特征比较带来一定的困难;基于特征空间的分析方法主要是基于统计规律得出判断,因此受噪声的影响相对较小,但却需要一定数量的训练样本和先验知识。

基于特征空间的高光谱数据分析方法已有许多年的研究基础,尤其是对多光谱数据,有许多成熟和经典的分析方法,包括 k 近邻法^[22]、贝叶斯 (Bayes) 分类器^[23]、高斯过程分类器^[24]、神经网络分类器^[25]以及支持向量机 (support vector machine, SVM) 分类器^[26]等,相应的特征提取算法也有主成分分析法 (principal component analysis, PCA)^[27]、线性判别分析法 (linear discriminant analysis, LDA)^[28]等经典方法。但对高光谱数据而言,至今还未出现有针对性的、具有更好效果的基于特征空间的数据分析方法。专业高光谱数据处理软件,如 ENVI、ERDAS 所提供的基于特征空间的分析方法也只有多光谱数据中常用的几种方法。由于高光谱数据有不同于多光谱数据的特点,如何让多光谱数据基于特征空间的分析方法在高光谱数据中应用是一个值得研究的问题。

1.1.2 高光谱数据应用领域

与多光谱数据技术相比,高光谱数据技术在各个领域都具有巨大应用潜力,主要体现在以下几个方面^[29,30]。

1. 测绘和更新地形图

高光谱数据能够获得精细的地物光谱,可以广泛地应用于绘制和更新地形图。如:矿物勘查、矿物成分识别、岩石识别等地质岩矿的分布图,农作物、森林、草地等植被精细分布图,城市建设的变化图等。

2. 环境、食品与灾害的快速监测

不同成分的物质反射的光谱特征是不同的,进行环境安全、食品卫生与自然灾害的快速监测是促进国家建设发展中不可缺少的环节。如:污染物分布调查,水质与大气污染物监测,城市绿化带覆盖度调查,食品生产过程中卫生控制,地震、洪

涝、虫害、火灾、沙尘暴等自然灾害快速监测及预防,城市热岛效应监测等。

3. 大气与海洋遥感应用

通过高光谱数据探测出大气与海洋中不同物质成分的反射和吸收规律,并反映在光谱曲线中。所以在探测大气与海洋的主要成分基础上,可以进行天气预报,台风预测,海洋资源勘查,水色、水温变化监测,海岸带变化、海洋生态及污染监测等。近年来典型的应用,如高光谱数据在太湖蓝藻监测中的应用^[31]。

4. 军事领域应用

利用高光谱数据可以揭露伪装、隐藏和欺骗的军事目标,毁伤效果分析,获取其他感兴趣的军事情报等。

5. 行星探索

通过航天器搭载的成像光谱仪,对外星球上资源信息、生命迹象、水含量等进行探索。

1.1.3 高光谱数据降维的研究现状及存在问题

高光谱数据量很大,包含不同地物的信息非常丰富,但其谱间相关性强,且也存在很多冗余信息。这些冗余信息不仅会影响地物识别和分类的效果,而且还会增加数据处理的代价。Chang 发现在不影响整体分类精度的前提下,最高可能有 94% 的光谱波段是没有必要的^[7]。因此,在实现地物识别和分类前,首先需要对接高光谱数据进行降维处理,以便保留有“价值”(便于进行高光谱数据分类识别)的地物信息,减少冗余信息,提高地物识别和分类的效率。一般来说,降维算法主要从两个方面进行:特征选择是直接从原始波段空间选择若干波段用于后续处理,约简后的特征集是原始数据的子集;特征提取则是对原始数据集中一个或若干个原始波段按照一定的操作函数进行变换,然后选择若干分量作为后续处理应用的特征子集。

一些传统的特征选择算法包括基于信息熵(联合熵)的选择^[32,33],基于分形维数的最佳波段指数选择^[34],基于波段相关度、离散度或巴氏距离等的选择^[35]等。这些方法往往试图对所有波段选择最优组合,但研究表明,以最佳波段指数、联合信息熵等对全部波段进行搜索计算的最优搜索方法在高光谱数据中因为计算量太大的原因难以得到应用,因此往往要研究次优选择算法。最常用的次优选择算法有顺序前向选择法^[36]、顺序后向选择法^[37]和最速上升搜索算法^[38]。随着计算智能、进化计算等理论的发展,粗糙集^[39]、遗传算法^[40]和蚁群优化算法^[41]等新方法在高光谱数据的降维处理中也陆续得到了应用。但是,由于特征选择受搜索方法和决策准则的显著影响,无论如何选择都必然会损失大量信息,因此更多的研究工作倾

向于特征提取。

通过特征提取技术,原始高维高光谱数据被映射或变换至低维空间(同时仍保留原始数据的某些必要特征),从而可在很大程度上避免维数灾难,使后续分类或聚类任务不仅更加稳定、高效、易于处理,而且更为重要的是,产生更优的泛化性能。目前,已有众多特征提取方法先后被提出并应用于高光谱数据的降维,如最小噪声分离(minimum noise fraction, MNF)^[42]、投影寻踪^[43]、小波变换^[44]、主成分分析(PCA)^[45]、线性判别分析^[46]、独立成分分析(independent component analysis, ICA)^[47]等。这些方法具有坚实的理论基础,易于执行和分析,得到了许多成功的应用,如 PCA 和 MNF 方法已成为一些商业化软件的标准操作模块。但是,它们均为(全局的)线性方法,无法揭示数据内在的非线性结构,而高光谱数据是本质非线性的。为了实现高光谱数据的非线性特征提取,可以借助于核技术,将传统的线性技术核化,如 Yang 等提出的核 Fisher 判别分析^[48]、Fauvel 等提出的核 PCA^[49]以及 Bai 等提出的核 ICA^[50]。但是,核化的特征提取方法往往依赖于某种隐式映射,不易直观地理解其工作机理,并且如何选择核及配置最优的核参数尚无可靠的理论依据。另一类重要的非线性特征提取技术是基于局部特性的流形学习方法,如罗琴等提出的局部线性嵌入^[51]和 Wang 等提出的等距特征映射 ISOMAP^[52]。为克服流形学习通常存在的“out-of-sample”问题^[53],即对样本集之外的数据点必须重复执行原有的整个学习过程或者求助于特殊的处理技巧,Chen 等给出了基于局部保持投影(locality preserving projection, LPP)的高光谱数据特征提取算法^[54]。LPP 本质上是拉普拉斯特征映射的线性化版本,既具有线性方法简单、快捷、可延展的优点,又具有一般线性方法所不具备的非线性流形学习能力,在高光谱数据特征提取领域得到了较好的应用。但是,LPP 需要付出参数选择的额外代价,并且最近的研究表明,参数的微小变化将导致最终结果大相径庭^[55]。虽然交叉验证是常用的参数选择技术,但往往只适合于监督学习,并且耗费大量训练样本,导致较高的计算开销。事实上,当训练样本(特别是有标记训练样本)较少时,目前尚无可靠的方法进行参数选择。

综合分析,现有的高光谱数据降维技术主要存在下面几个问题。

1. 如何更好地进行非线性特征提取

如上所述,现有的高光谱数据非线性特征提取方法存在“out-of-sample”问题、参数选择的额外代价及计算过于复杂。因此,有必要研究无参数或参数少和计算量较少的非线性特征提取方法。

2. 如何针对高光谱数据的三维特性设计降维算法

上文提到的降维算法,主要是将高光谱数据处理成向量的形式后再进行降维的。而高光谱数据是三维的,在处理成向量的过程中,难免会有信息的损失。故需

要设计一种直接对三维的高光谱数据进行降维的张量降维算法。

3. 如何同时处理呈现出“三全”“三高”“三多”发展趋势的大量高光谱数据

如果对这些高光谱数据进行监督降维,就需要对每个高光谱数据收集足够数量的标记训练样本,这不仅费时耗力,而且成本很高,是不现实的。因此,需要设计一种依靠单次采集的数据来处理一系列数据的降维算法。

4. 在标记样本很少的情况下,如何较好地选择训练样本和提取地物的本质特征

随着数据采集技术和存储技术的发展,获取无标记样本已变得非常容易。另一方面,由于有标记样本的获取需要相关领域的专家对样本进行标记,因而相对比较困难而且代价昂贵。所以在实际的高光谱数据降维应用中,通常会有大量的无标记样本,而有标记样本只占很小的比例。如何充分利用有限的标记样本和大量无标记样本来提高高光谱数据的降维效果,也是值得探讨的问题。

1.2 降维算法的研究现状

前面已经针对经典的降维算法进行了简介,本节重点介绍三种较为先进的降维算法研究,也是本书主要研究的对象:基于稀疏表示的降维算法,张量型降维算法,基于特征的迁移学习方法。

1.2.1 基于稀疏表示的降维算法

稀疏表示 (sparse representation, SR) 是近年来信号处理和模式识别领域的一个研究热点,是对多维数据进行线性分解的一种表示方法^[56]。它的稀疏性表现在对每个输入的信号,只有少数几个基函数具有较大的响应输出,而其他基函数的输出接近于零。因此,稀疏表示在图像降噪^[57]、修复^[58]、超分辨率处理^[59]、压缩感知等经典的图像和信号处理问题上表现出了优越的性能。近年来,随着机器学习和模式识别领域的发展,考虑到稀疏表示具有自然的判别能力^[56],能获得相互独立的特征,同时系数的稀疏分布能更好地拉开各类特征之间的距离,稀疏表示被推广到降维^[60]、分类^[61]、目标探测^[62]等相关领域。另外,降维算法的主要目的是:在保证一定学习性能的前提下提取尽可能少的特征数目。因此,降维算法在某种意义上来说也是一种稀疏学习方式。

目前,求解稀疏表示方法主要有 lasso^[63]、lars^[64]、elastic net^[65]。Zou 等在原始 PCA 上引入 lasso 和 elastic net 稀疏方法,提出稀疏主成分分析 (sparse PCA)^[66]。类似的, Clemmensen 等提出稀疏判别分析^[67], Qiao 等提出稀疏线性判别分析^[68], Zheng 提出稀疏局部保持嵌入^[69]。而 Moghaddam 等将谱边界和稀

疏子空间学习融合在一个框架中,即利用贪婪算法和广义谱边界的系数主成分分析^[70]和稀疏线性判别分析^[71]两种算法。随之,Cai等将谱回归方法融入到经典的子空间学习中,如PCA、LDA和LPP,提出一种新的降维框架统一稀疏子空间学习方法^[72]。2000年,Cai等^[73]在AAAI国际会议上发表了基于图的稀疏投影方法。在此基础上,Lai等提出在保持稀疏关系的同时最大化不同样本间距离的稀疏局部判别投影^[74]。Zhou等提出流形弹性网络(manifold elastic net, MEN)方法,并利用MEN提出稀疏降维的一种框架^[75]。Wright等^[76]、Cheng等^[77]、Huang等^[78]和Qiao等^[79]研究员先后利用稀疏表示构建 l_1 图,并应用到子空间学习方法中,即寻找一个能保留原始高维数据稀疏关系的低维子空间,此处统称稀疏保持投影(sparsity preserving projections, SPP)。随后,Qiao等将半监督判别分析(semi-supervised discriminant analysis, SDA)^[80]中正则项用SPP代替,提出稀疏保持判别分析(sparsity preserving discriminant analysis, SPDA)的半监督方法^[81]。Wong在SPP基础上引入非负矩阵分解,提出具有自然判别信息的非负稀疏保持嵌入(non-negative sparseness preserving embedding, NSPE)^[82]。Gui等^[83]和Lu等^[84]在SPP基础上引入判别信息,分别提出判别稀疏近邻保持嵌入两种监督方法。

鉴于上面的方法在人脸识别中表现突出,本书取其精华提出新的非负稀疏子空间学习算法^[85-87],并推广应用到高光谱数据降维处理中。

1.2.2 张量型降维算法

实际上,高光谱数据是立方体结构的,而1.2.1节中所提到的降维算法都是将数据转换成向量形式处理的,这样会破坏原始数据的空间关系。为此,研究者们开始研究张量型降维算法。此处主要介绍三阶以上的张量子空间分析方法。

目前,张量子空间分析方法主要有:Lu等将PCA推广到任意高阶张量空间中,提出多线性主成分分析(multilinear PCA, MPCA)^[88];Yan等将带有判别信息的LDA方法推广到高阶张量空间中,提出多线性局部判别分析(multilinear discriminant analysis, MDA)^[89];Tao等以最大间距准则(maximum margin criterion, MMC)^[90]作为准则函数,并加入调节函数,提出广义张量判别分析^[91];Zhang等^[92]将补丁校准框架推广到高阶张量空间中,提出张量判别局部校准算法。

1.2.3 基于特征的迁移学习方法

前面提出的方法只能解决相同数据集中的问题,而在现实世界中,往往存在大量的不同分布、不同形式的数据。为了同时处理这种类型的数据,在机器学习中,可以通过领域自适应^[93]和迁移学习^[94]来解决。迁移学习的目的是解决当来自一个或多个源领域的训练样本和来自目标领域的测试样本属于不同分布或是不同特征空间表示时的问题。迁移学习的关键思想是:虽然源和目标领域之间的分布不