

面向高寒湿地分类的遥感特征优选研究

霍轩琳^{1, 2}, 牛振国², 张波², 刘林崧², 李霞³

1. 长安大学地球科学与资源学院, 陕西 710054;

2. 中国科学院空天信息创新研究院 遥感科学国家重点实验室, 北京 100094;

3. 长安大学土地工程学院, 陕西 710054

摘要: 高寒湿地是青藏高原重要的地表覆盖类型之一, 对于水源涵养、调节气候、维护生物多样性等起着关键作用, 准确及时获知高寒湿地的时空分布对于湿地的保护和管理十分必要。遥感分类特征优选对湿地制图具有关键性的作用。虽然像光谱特征、纹理特征、地形特征等均在已有研究中有涉及, 但鲜有研究聚焦光谱指数特征, 深入探讨其数理统计特征和特征优选方法。本研究以甘肃首曲高寒湿地保护区为研究区, 基于 Sentinel-2 数据得到各分类特征(光谱、植被指数、红边指数和水体指数), 采用 Filter 和 Wrapper 特征选择方法包括 Jeffries-Matusita 距离、光谱角距离(SAD)、欧氏距离(ED)、RF-RFE 算法和 Relief-F 算法对上述特征进行优选, 并利用 Filter 方法的 Z 检验进行量化评价。研究表明:(1) 所有参与分类的类别中, 河流与裸地最容易区分, 其次为草原与沼泽, 沼泽化草甸与草甸较为难分。对沼泽、沼泽化草甸、草甸、草原邻近两类可尝试 MCARI2、NDWI、DVI、EVI、EWI、IRECI、MCARI、TCARI、UGWI 指数进行区分;(2) 就不同指数特征对湿地信息提取贡献程度而言, 水体指数特征>植被指数特征>红边指数特征;(3) 从特征优选方法角度看, Filter 方法中的 ED 距离算法与 Relief-F 算法表现突出;(4) 最终选出适于高寒湿地信息提取的指数有 RDVI、NDVI、MSR、RVI、VIgreen、RNDWI、NDWI、NDWI_B、MNDWI、EWI、CIre;(5) 从不同分类特征的数理统计指标看, 中值特征的分类结果最好, 其次是平均值特征。本研究为湿地信息提取在特征变量优选方面提供了一种可迁移且普适性高的方法和思路。

关键词: 湿地分类; 高寒湿地; 特征优选; 青藏高原; Sentinel

中图分类号: 文献标志码: A

引用格式:

[DOI:]

收稿日期: ; 预印本:

基金项目: 第二次青藏高原综合科学考察研究项目(2019QZKK0106); 基于时间序列遥感数据的中国湿地分类方法研究(41971390) 资助。

第一作者简介: 霍轩琳, 1998 年生, 女, 硕士研究生, 研究方向为湿地生态与环境遥感研究。E-mail: h1160964174@gmail.com

通信作者简介: 牛振国, 研究员。E-mail: Niuzg@aircas.ac.cn

1 引言

湿地具有调节气候、保护生物多样性、蓄洪抗旱、改善环境等功能，同时为动植物提供了良好的生存条件（Moor 等，2015），是地球上最重要的生态系统（森林、海洋、湿地）之一（何菊红等，2015）。青藏高原高寒湿地是青藏高原乃至西部地区最重要的生态系统，近年来受到自然和人为因素干扰，高寒湿地面积已锐减了 10%，且水量和湿地的面积减少速度还在加快（王根绪等，2007；徐新良等，2008）。因此，及时获知高寒湿地面积、分布区域等对青藏高原高寒湿地管理与保护乃至生态系统的可持续发展至关重要。

遥感分类特征的选择是目前对湿地进行准确分类与制图的众多挑战之一。分类特征需要考虑到以下两方面，其一，湿地是陆地生态系统和水域生态系统的交界地带，具有较高的景观异质性，仅靠一种特征变量可能无法很好地进行湿地的准确提取；其二，若使用过多的特征变量参与其中将影响分类精度和效率。可见，多特征变量提取与优化以及进行有效组合将是今后湿地信息智能化提取的重点难点（张磊等，2019）。特征选择（Feature Selection）通常情况下是将特征按照相关性准则排序，去掉冗余和不相关的特征（Guyon 等，2003），按评价标准不同，特征选择算法可分为过滤式（Filter）、封装式（Wrapper）和嵌入式（Embedded）三种（M. Dash 等，1997；M. Dash 等，2002；Saeys Y 等，

2007；Kira K 等，1992）。其中，Filter 方法利用特征本身内在特性给出特征评价，特征评价越高表示该特征区分能力越强。这也是该方法最主要的特点：特征选择独立于分类学习算法。它不依赖某种分类器，因此简单，速度快，效率高。L.Yu 等（2003）以特征间的相关性为指导，用 Filter 方法进行特征选择，证实了该方法进行特征选择的有效性；Wrapper 方法是将特定学习算法性能作为筛选子集的评估准则，每次筛选出的特征子集都需调用特定分类器进行精度验证。G John 等（1994）提出当满足一定条件时，将获得具有较高分类性能的识别模型。该方法准确率较高，利于关键特征的识别，在算法速度上比 Filter 方法慢，时间复杂度较高；Embedded 方法是一种基本的归纳方法，可以说是 Wrapper 方法的延伸。Embedded 方法将特征选择过程嵌入到分类器的建造过程中，主要的例子是套索回归的问题以及决策树如 Breiman 的 CART 算法（Breiman，2001）。该方法计算效率高，但是特征中可能存在无关特征降低分类精度，在本次研究中未涉及。

针对已有的特征选择方法，有一些学者就其在湿地遥感分类方面进行了探索。如 Masoud Mahdianpari 等（2019）利用 JM 距离定量地确定不同类型湿地的可分性，结合随机森林进行分类总体精度达到 88.37%；孙艳丽等（2015）利用光谱角距离（SAD）和欧氏距离（ED）双重判定提取不变特征点，提出了一种基于光谱角—欧氏距离的辐射归

一化方法；郝玉峰等（2021）利用 Relief-F 算法计算了 52 个特征变量的权重，选出前 20 个特征变量构成最优特征集参与湿地信息提取；Han 等（2011）利用 Z 检验方法测定区分两种植被类型的最佳纹理波段。以上提到的这些方法均属于 Filter 方法；解淑毓等（2021）采用 Wrapper 方法中典型的 RFE 算法进行沼泽湿地分类中的变量优选，显著减少了数据冗余。Phan 等（2020）在其研究中指出 GEE 提供的像元级重组规则主要包括最大值、最小值、平均值、中位数和百分位数等。

总体而言，目前湿地分类特征优选的研究多集中于通过多特征变量参与、单一特征优选方法来甄选最优特征集。不同特征的统计方式和不同特征优选方法对分类的影响尚未见相关研究报道，同时不同分类特征对高寒湿地类别分类的适用性也未见相关研究。鉴于此，本文基于 Sentinel-2 影像数据，以首曲高寒湿地保护区为研究区域，利用随机森林分类算法，探讨数理统计特征和特征优选方法对优选的各种影响，并分析不同特征对高寒湿地类型分类的适用性。该研究将对提高高寒湿地遥感制图具有重要的参考价值。

2 研究区概况与数据源

2.1 研究区概况

甘肃黄河首曲国家级自然保护区（33° 20′ 01″ N—33° 56′ 31″ N，101° 54′ 12″ E—102°

28′ 45″ E）位于甘南藏族自治州玛曲县境内，属于内陆湿地和水域生态系统类型的自然保护区，是青藏高原典型的面积较大的高寒湿地（薛鹏飞等，2021），也是全球保存状态最为完整和原始的湿地。首曲高寒湿地属于高原大陆性气候，年均气温 1.1℃，年平均降水量 615.5mm，全年降雨 150 天左右（高斌斌，2008），给黄河贡献了黄河源区总径流量 58.7% 的水量，被称为“蓄水池”和“高原水塔”。

2.2 数据及预处理

2.2.1 数据源及预处理

近年来随着遥感技术的迅猛发展，越来越多的多源传感器涌现，其空间分辨率、时间分辨率、波段数量得到了巨大提升（郑阳等，2017），为湿地的遥感分类提供了更多的选择。综合考虑影像分辨率、波段、可获得性等多因素，研究采用 Sentinel-2 影像。Sentinel-2 属于中高空间分辨率遥感影像，携带高分辨率多光谱传感器 MSI，可提供可见光、近红外到短波红外的 13 个波段，是目前唯一在植被光谱的红边区域（670-760nm）设置 3 个波段的卫星。红边波段数据及其衍生指数可区分 C3、C4 植被（Shoko 等，2017；Korhonen 等，2017；CHANG 等，2020），极大促进了对植被生长信息及其健康状况的有效监测（张磊等，2019）。本研究使用的 Sentinel-2 数据是用 GEE 平台“COPERNICUS/S2_SR”数据集中 2020 年 1 月 1

日至 2020 年 12 月 31 日的影像数据，采用已经过辐射定标和几何校正的 Level-1C 产品，去云层覆盖率大于 10% 的影像后得到 34 景 Sentinel-2 的无云影像。

2.2.2 样本数据

本研究以 Global Lakes and Wetlands Database、Wetland Dataset of CAS 湿地制图产品公开数据集为参考数据集，在 Google Earth 软件上开展样本集目视解译工作。最终取得样本点共 480 个，每个地类（分类体系见表 1）各 80 个(图 1)。

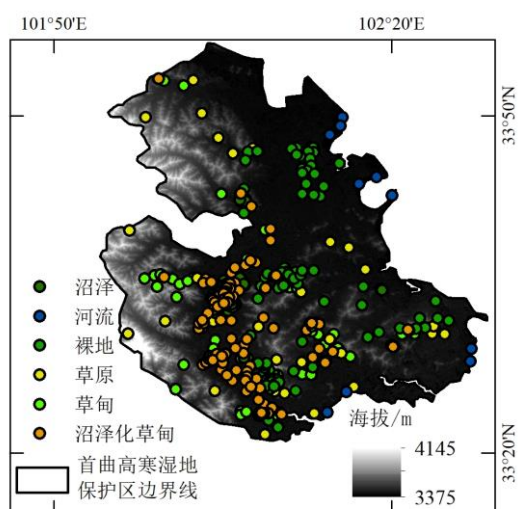


图 1 黄河首曲国家级自然保护区样本点分布图

表 1 黄河首曲自然保护区土地覆被分类方案及影像示例


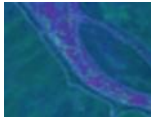
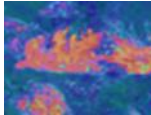
Table 1 Land Cover Classification Scheme and Wetland Corresponding Image of Yellow River Shouqu Nature Reserve

分类	具体说明	Sentinel-2 影像举例
沼泽	地表及地表下层土壤经常过度湿润，地表生长着湿性植物和沼泽植物，有泥炭累积或虽无泥炭累积但有潜育层存在的土地。	
沼泽化草甸	湿中生多年生草本植物为主的植物群落，为典型草甸向沼泽植被的过渡类型，土壤为暗色草甸土。	
草甸	向草原过渡地带的典型草甸，可视为地带性植被。土壤为排水良好，湿度中等的草甸黑土。	

Fig.1 Distribution map of sample points in the Yellow River Shouqu National Nature Reserve

2.3 分类体系方案

Lori White 等（2020）指出，草甸沼泽对水位变化特别敏感，下垫面含水量一定程度上决定了湿地植被构成。本次不直接以湿地所在地域或湿地中细分植被为分类标准而是优先考虑造成湿地植被生长的起因——下垫面土壤含水量。随着土壤水分的增减，草原地区的草甸可能会发生演变，当水分增加时，可转变为沼泽；当水分减少时可转变为草原，沼泽化草甸是草甸与沼泽之间的过渡类型。沼泽、沼泽化草甸、草甸区域及边界的变化是监测湿地变化的重要指标，其变化可反映出当地的自然气候，生态环境的变化。因此本文土地覆被分类方案如表 1 所示。

草原	以多年生旱生型草本植物占优势的水平或垂直地带性植被类型。	
水体	主要包括湖泊：陆地表面洼地积水形成的比较宽广的水域；河流：指地表上有相当大量且常年或季节性流动的天然水流。	
裸地	地表有土壤覆盖且植被覆盖度小于 5% 的土地。	

3 研究技术流程和方法

3.1 技术流程

本文的技术路线图如下图 2 所示。首先对 Sentinel-2 影像数据进行预处理。基于样本计算 32 种特征指数（3.2 部分）的统计特征（均值、标准

差、中值、最大值、最小值）后，分别利用 JM 距离、ED 距离、SAD 距离、RF-RFE 算法、Relief-F 算法进行遴选，获取不同特征优选方法下的最优特征集。基于不同优选特征集利用随机森林进行湿地分类，依据特征优选方法及随机森林分类结果评价不同特征优选方案。

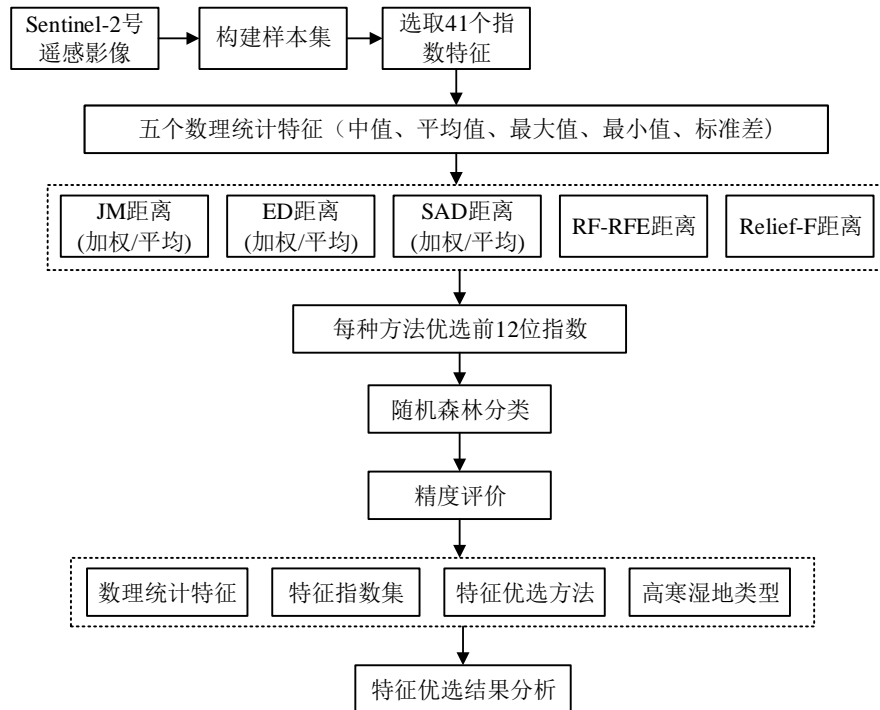


图 2 特征优选技术流程图

Fig.2 Feature selection technology flow chart

3.2 特征说明

本文选择了能够表征湿地植被、水文和土壤特征的特征、植被指数、水体指数、红边指数（表 2）。的指数进行分析，共选取了 32 种特征集，包括光谱

特征、植被指数、水体指数、红边指数（表 2）。

表 2 Sentinel-2 特征集概述

Table 2 Sentinel-2 feature sets list

特征	指数	Sentinel-2A 计算公式	参考文献	指数描述
变量	简称			
光谱	Band	B2、B3、B4、B5、B6、B7、B8、B11、B12		
特征				
	NDVI	$(B8 - B4)/(B8 + B4)$	Defries 等, 1994	对绿色植被的生长状况具有良好的表征能力, 表征植被覆盖度和生长与健康状况。
	VIgreen	$(B3 - B4)/(B3 + B4)$	Gitelson 等, 2002	通过冠层背景调节、气溶胶阻力系数和增益因子补偿土壤背景和大气效应, 且对高生物量更敏感。
	EVI	$2.5 \times \left(\frac{B8 - B4}{B8 + 6 \times B4 - 7.5 \times B2 + 1} \right)$	Huete 等, 2002	RDVI 取 DVI 和 NDVI 两者之长, 可用于高低不同植被覆盖的情况。
	RDVI	$(B8 - B4)/(\sqrt{B8 + B4})$	Roujean 等, 1995	该指数根据对两个光谱带组合得出的几种植被指数的评估而制定, 包括 NDVI、SR、SAVI、SAVI1、SAVI2、WDVI、GEMI、NLI 和 RDVI。
	MSR	$\left(\frac{B8}{B4} - 1 \right) / \left(\sqrt{\frac{B8}{B4} + 1} \right)$	Chen 等, 1996	对叶绿素含量变化非常敏感。
植被	TCARI	$3 \times \left[(B8 - B4) - 0.2 \times (B8 - B3) \times \left(\frac{B8}{B4} \right) \right]$	Haboudane 等, 2002	该指数基于叶绿素吸收导致红色反射率降低, 叶组织丰度导致 NIR 反射率增加, 这两者都会增加三角形的总面积。
指数	TVI	$0.5 \times [120 \times (B8 - B3) - 200 \times (B4 - B3)]$	Broge 等, 2001	用于估算和检测植被生物量, 并对高植被覆盖度反应敏锐。
	RVI	$B8/B4$	Birth 等, 1968	DVI 对土壤背景的变化较 RVI 敏感, 植被覆盖度高时, 对植被的灵敏度有所下降。
	DVI	$B8 - B4$	Tucker 等, 1979	GCVI 比 NDVI 具有更大的动态范围, 适合高密度植被覆盖区域。
	GCVI	$\frac{B8}{B3} - 1$	Gitelson 等, 2003	消除不同背景的土壤所反射的不同光谱对植被特性的影响, SAVI 包含了土壤调节系数, 更适用于低植被覆盖地区。
	SAVI	$1.5 \times (B8 - B4)/(B8 + B4 + 0.5)$	Huete 等, 1988	

	gNDVI	$(B8 - B3)/(B8 + B3)$	Gitelson 等, 1996	gNDVI 与叶绿素含量和叶面积指数有显著相关性。
	MCARI	$[(B8 - B4) - 0.2 \times (B8 - B3)] \times (B8/B4)$	Daughtry 等, 2000	对叶片叶绿素浓度和背景反射率都有响应。
	GI	$B3/B4$		
	NDWI	$(B3 - B8)/(B3 + B8)$	McFeeters 等, 1996	能够抑制植被信息, 突出水体; 对建筑物和土壤的分离有一定影响; 受冰雪、薄云和山体阴影影响较大, 适用于地势平坦地区。
	NDWI_B	$(B2 - B4)/(B2 + B4)$	ZHANG 等, 2020	可以在一定程度上抑制与水体无关的背景信息, 增强水体信息。
	MNDWI	$(B3 - B11)/(B3 + B11)$	Xu 等, 2005	能够较好去除居民地和土壤等影响, 突出水体; 受阴影影响大。
水体 指数	RNDWI	$(B12 - B4)/(B12 + B4)$		能削弱混合像元和山体阴影的影响, 较好地提取水陆边界; 一般适用于山区等地形起伏较大地区。
	LSWI	$(B8 - B11)/(B8 + B11)$	Xiao 等, 2005	对植被冠层含水量和土壤含水量敏感。
	EWI	$(B3 - B8 - B12)/(B3 + B8 + B12)$	Yan 等, 2007	能够抑制居民地、土壤和植被等噪声; 易受到阴影及浅滩的影响; 适合半干旱地区的水体提取。
	UGWI	$\frac{B3^3 - B2 \times B4 \times B8}{B3^3 + B2 \times B4 \times B8}$	DUAN 等, 2021	能较完整地提取不同泥沙含量的洪水、坑塘水面及阴影下各种水体, 不易受阴影和暗色地物等干扰。
	SWI	$B2 + B4 - B8$	Chen 等, 2015	较好地区分水体和阴影, 能削弱积雪和山体裸地的影响, 适用于山区的水体提取。
	NDVIre1	$(B8 - B5)/(B8 + B5)$	Gitelson 等, 1994	用波段 5 替代 NDVI 中的红波段。
红边 指数	NDVIre2	$(B8 - B6)/(B8 + B6)$	Gitelson 等, 1994	用波段 6 替代 NDVI 中的红波段。
	NDVIre3	$(B8 - B7)/(B8 + B7)$	Gitelson 等, 1994	用波段 7 替代 NDVI 中的红波段。
	NDre1	$(B6 - B5)/(B6 + B5)$	Gitelson 等, 1994;	用红边的峰和谷来代替传统 NDVI 中的红

NDre2	$(B7 - B5)/(B7 + B5)$	Barnes 等, 2000	光和近红外波段, 可用于估算植物叶面积指数和叶绿素含量。
CIre	$\frac{B7}{B5} - 1$	Gitelson 等, 1994;	将 NDVI 中的近红外和红波段替换为波段 7 和波段 5, 它可用于精细农业、森林监测、植被胁迫性探测等。
MTCI	$(B6 - B5)/(B5 - B4)$	Barnes 等, 2000	该指数与植物叶绿素, 氮素含量具有显著的线性关系。
MCARI2	$((B6 - B5) - 0.2 \times (B6 - B3)) \times \left(\frac{B6}{B5}\right)$	Dash 等, 2004	对植物叶片叶绿素含量较敏感, 其值越大代表叶绿素含量越高。
REP	$705 + 35 \times \frac{0.5 \times (B4 + B7) - B5}{B6 - B5}$	Wu 等, 2008	该指数对植物中的叶绿素含量较为敏感, 其值越大表示叶绿素含量越高。
IRECI	$(B7 - B4)/(B5/B6)$	Guyot 等, 1988	当植物叶片的叶绿素含量增加时, REP 向长波方向移动, 反之则向短波方向移动。
		Frampton 等, 2013	该指数与植物冠层叶绿素含量和叶面积指数具有相关关系, 可定量表征植物的叶绿素含量。

3.3 特征优选方法

本文选择包括 JM 距离、ED 距离、SAD 距离、Relief-F 算法、Z 检验在内的五种 Filter 方法, 以及 Wrapper 方法的 RFE 算法进行本次实验。

3.3.1 Jeffries-Matusita 距离

Jeffries-Matusita (JM) 距离基于数据正态分布的假设得到不同类别的分离度, 在模式识别和特征选择领域中较为广泛的使用 (Dabboor 等, 2014)。对训练样本集 $C(i, j = 1, 2, \dots, C, i \neq j)$ 中两个待分地类 w_i 和 w_j 之间的 JM 可分性准则定义如下:

$$J_{ij} = 2(1 - e^{-d_{ij}}) \quad (1)$$

其中, d_{ij} 是两个待分地类 w_i 和 w_j 之间的 Bhattacharyya 距离, 定义为:

$$d_{ij} = -\ln\left\{\int \sqrt{P(x/w_i)P(x/w_j)} dx\right\} \quad (2)$$

上式中, $P(x/w_i)$ 和 $P(x/w_j)$ 是地类 w_i 和 w_j 的随机变量 x 的条件概率密度函数, 通常假设多元正态分布, Bhattacharyya 距离公式表达为:

$$d_{ij} = \frac{1}{8}(m_j - m_i)^T \left[\frac{\Sigma_i + \Sigma_j}{2}\right]^{-1} (m_j - m_i) + \frac{1}{2} \ln \frac{|\Sigma_i + \Sigma_j|}{\sqrt{|\Sigma_i||\Sigma_j|}} \quad (3)$$

其中, m_i 和 m_j 分别表示均值; Σ_i 和 Σ_j 分别表示 w_i 和 w_j 的协方差矩阵, 上标 T 表示矩阵的转置。

3.3.2 欧氏距离

欧氏距离 (Euclidian distance, ED) 是常见的

相似性度量方法，其实质是通过一定的准则函数，求两个不同地类的像元对应的光谱向量之间的距离，此距离代表两像元的差异程度。两种地类的欧式距离越大，代表两种待分地类间的可分性越强，反之，则表示可分性越弱（Osmar 等，2011）。由于欧式距离算法默认每一个维度是相同权重，因此如果不同维度取值范围差别较大时需要先对其进行归一化，ED 值计算公式为：

$$ED = \sqrt{\sum_{i=1}^N (X_i - Y_i)^2} \quad (4)$$

式中：i 表示波段，N 表示波段总数， X_i 和 Y_i 分别表示两种待分地类样本集所对应的像元亮度值。

3.3.3 光谱角距离

光谱角距离（Spectral Angle Distance, SAD）是常用的光谱分类方法。在光谱空间中，每个像元对应一个多维光谱向量，将两个向量之间的夹角定义为光谱角。光谱角越小，两光谱越相似，属于同类地物的可能性越大。由于光谱角距离不受光照、阴影等条件的影响，因此可以突出目标光谱形状特征。两种待分地类光谱的相似度越高，SAD 值越大，最大取值为 1(Kruse 等，1993)。计算公式为：

$$\theta = \cos^{-1} \frac{\sum_{i=1}^N X_i Y_i}{\sqrt{\sum_{i=1}^N X_i^2 \sum_{i=1}^N Y_i^2}} \quad (5)$$

$$SAD = \cos \theta \quad (6)$$

式中：i 表示波段，N 表示波段总数， X_i 和 Y_i 分别表示两种待分地类样本集所对应的像元亮度值。

3.3.4 Z 检验方法

Z 检验方法能够用来测定两种地物类型在不同特征变量间统计显著性差异。具体步骤为首先将特征变量分为植被指数、水体指数、红边指数三组，再分别计算两种湿地类型在不同变量的 Z 值。Z 统计表达式如下(Ning 等，2011)：

$$Z = \frac{(u_1 - u_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (7)$$

其中， u_1 和 u_2 指的是两种待分地类的平均像素值； n_1 和 n_2 指两种待分地类样本个数； s_1 和 s_2 指两种待分地类像素值的标准差。Z 值越大，待分地类在此特征上的差异就越显著。

以上四种方法均用于判别两地类之间的可分性，为直观地表达指数对所有地类区分的能力，进行以下步骤：根据在样本解译时发现的沼泽、沼泽化草甸、草甸、草原四种类型越相邻越难区分的认识，采用加权（表 3）对类型组合对应的 JM 值等进行处理，紧密相邻的两类权重赋 4，次相邻的两类权重赋 3，以此类推。在本次实验中同时采用了平均方法做补充。

表 3 各类型组合权重分配

Table 3 Weight distribution of various types of combinations

类型组合	权重分配	类型组合	权重分配	类型组合	权重分配
河流、裸地	1	沼泽化草甸、沼泽	4	草原、裸地	1
沼泽、裸地	1	草甸、裸地	1	草原、河	1

			流	
沼泽、河流	1	草甸、河流	1	草原、沼泽
				2
沼泽化草甸、裸地	1	草甸、沼泽	3	草原、沼泽化草甸
				3
沼泽化草甸、河流	1	草甸、沼泽化草甸	4	草原、草甸
				4

3.3.5 Relief-F 算法

Relief-F 算法是基于分析邻近样本对类别的区分能力继而确定特征的权重，核心思想是一个优秀的特征应该使得同类的样本更加靠近，而使得不同类的样本更加分散，它是 Relief 算法的拓展（刘吉超等，2021）。其原理为假设数据集 D 中有 N 个类别的样本，对属于第 n 类中样本 R ，首先在同类即第 n 类的样本中寻找 R 的 k 个最近邻样本 H ，作为猜中近邻；在第 n 类之外的每个类中均找到 R 的 k 个最近邻样本 M 作为猜错近邻，最后定义的权重为：

$$W(A) = W(A) - \sum_{j=1}^k \frac{\text{diff}(A, R, H_j)}{mk} + \sum_{C \neq \text{class}(R)} \left[\frac{p(C)}{1-p(\text{class}(R))} \sum_{j=1}^k \text{diff}(A, R, M_j(C)) \right] / (mk) \quad (8)$$

在上述公式中， $\text{diff}(A, R_1, R_2)$ 表示样本 R_1 和样本 R_2 在特征 A 上的差，其计算公式 $M_j(C)$ 表示类别 C 中的第 j 个最近邻样本， $p(C)$ 为该类别的比例。

3.3.6 递归特征消除法

相比递归特征消除法（Recursive feature elimination, RFE）（Elavarasan 等，2020），随机森林和 RFE 相结合形成 RF-RFE，其能够更加合理的决定最终特征子集的大小，避免了人为因素造成

的影响。RF-RFE 算法用于特征选择，是采用随机森林算法得到的重要性排序进行后向迭代删除特征重要度最小的特征，再将其余特征用随机森林算法重新评估后得到新的特征重要性排序，重复步骤，每次删掉特征重要性小的特征，最终得到分类的最优特征集。

3.4 分类方法

随机森林算法是由 Breiman 等人提出的一种统计学习理论（Breiman，2001），研究表明随机森林算法具有速度快，准确度高，稳定性好的优势。因此论文采用随机森林分类方法进行湿地信息分类。

$$H(x) = \arg \max_y \sum_k I(h_k(x) = y) \quad (9)$$

3.5 精度分析

混淆矩阵是一种特定的矩阵用来呈现算法性能的可视化效果，主要通过比较分类结果与实际测量值之间的混淆程度进行精度评价。本文利用混淆矩阵，分别计算总体精度（Overall Accuracy，OA）、Kappa 系数、生产者精度（Producer's Accuracy，PA）、用户精度（User's Accuracy，UA）。其中总体精度和 Kappa 系数作为评价总体分类精度的指标，生产者精度和用户精度作为衡量各类的漏分和错分误差的指标。

4 结果

为了方便叙述，论文用到的特征提取方法、数理

统计方法和数据处理方式分别按表 4 编码进行论述。

表 4 编号对照表

Table 4 Number comparison table

优选方法	代码	数理特征	代码	处理方式	代码
JM 距离	A	中值	a	加权	1
ED 距离	B	平均值	b	平均	2
SAD 距离	C	最大值	c		
RF-RFE 算法	D	最小值	d		
Relief-F 算法	E	标准差	e		

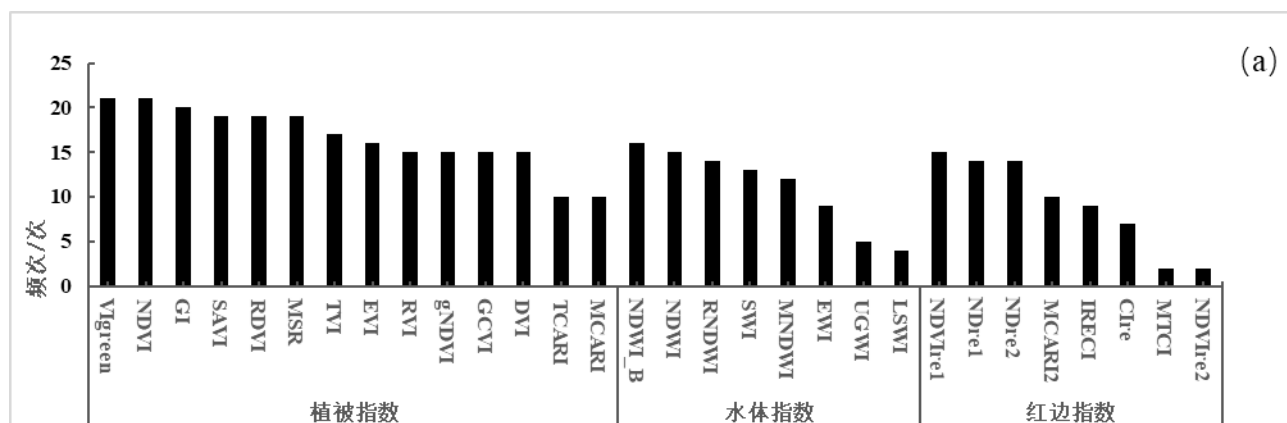
4.1 指数的优选

根据 JM 距离、SAD 距离、ED 距离、Relief-F 算法和 RF-RFE 算法计算结果，逐一得出在水体指

数、植被指数、红边指数中表现较好的指数。

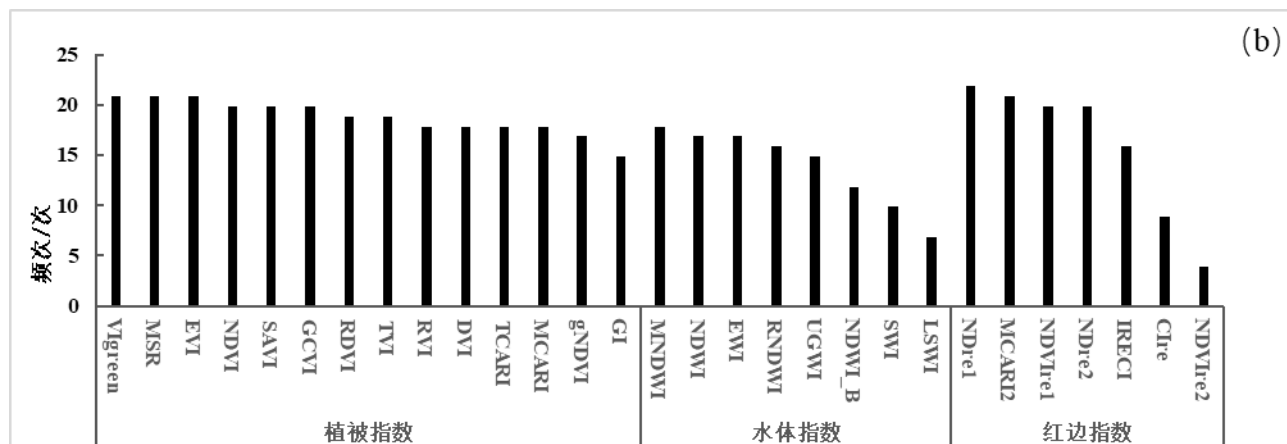
4.1.1 JM 距离

JM 距离取值范围为[0,2]，JM 距离大于 1.8 表示样本间可分性较好，据此选出每类组合中 JM 值大于 1.8 的指数（附表 1），并对这些指数进行频次统计（图 3-(a)）。植被指数在频次上大于水体指数和红边指数，其中最突出的是 NDVI、VIgreen 这两个指数；水体指数中频次数在前两位的指数为 NDWI_B、NDWI；红边指数中 NDVIre1、NDre1、NDre2 指数出现频次较多。



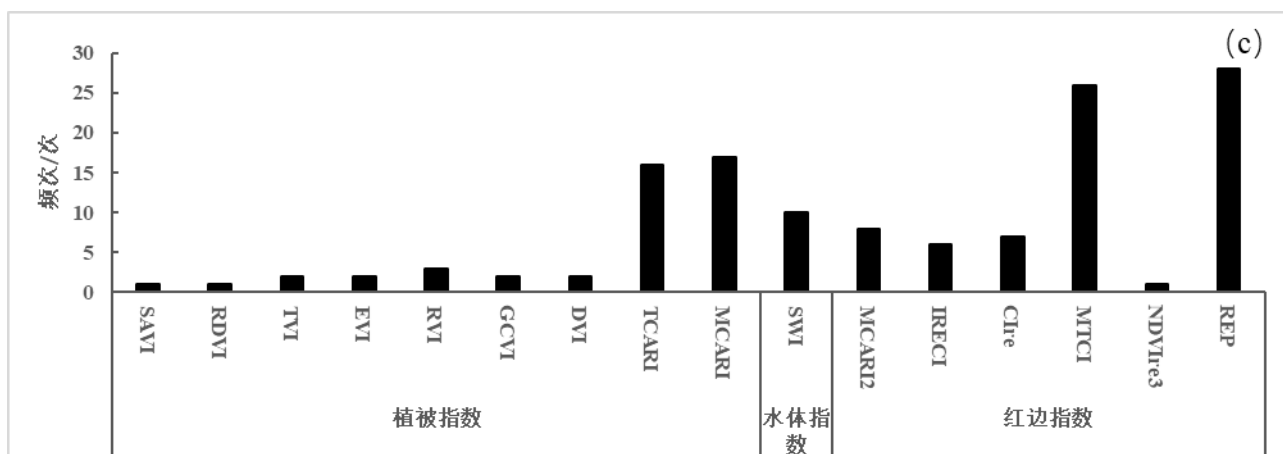
(a)JM 距离大于 1.8 的指数频次统计

(a)Exponential frequency statistics of JM distance greater than 1.8



(b)ED 距离大于 4 的指数频次统计

(b)Exponential frequency statistics of ED distance greater than 4



(c)SAD 距离小于 0.6 的指数频次统计

(c)Exponential frequency statistics of sad distance less than 0.6

图 3 JM 距离、ED 距离、SAD 距离指数频次统计

Fig.3 Frequency statistics of JM distance, ED distance and SAD distance index

采取加权和平均两种计算方式对数据进行处理，按重要性从大到小排序的前 12 位指数如附表 4 所示。发现加权方式优选的指数 TVI、DVI、SAVI 和以平均方式优选指数 RDVI、SAVI，不论在哪个数理统计特征中都有。当将两种方式所选的指数放在一起统计，优选至少出现五次的指数作为 JM 距离方法所选最优指数集，包括：SAVI、RDVI、TVI、DVI、MSR、EWI、NDWI、RNDWI、SWI、EVI、NDVI、gNDVI、RVI。基本包含了 JM 距离大于 1.8 所选频次较高的植被指数和水体指数，但是红边特征指数没能入选。

4.1.2 ED 距离

ED 值越大说明光谱距离越大，地类可分性越强。ED 值大于 4 表示样本间可分性好，依此标准对

优选的指数进行频次统计（附表 2），指数出现频次如图 3-(b)所示。红边指数 NDre1 出现次数最多。另外在红边指数中 MCARI2、NDre2、NDVIre1 均表现出较好的区分能力。植被指数 VIgreen、MSR、EVI 指数出现 21 次，其余指数与之相比相差不大，植被指数在分类中指数数量最多。水体指数 MNDWI、NDWI、EWI 表现优异。

同样对 15 种类型组合的 ED 值进行加权与平均两种方式计算，优选结果附表 5。加权方式优选的指数出现四次以上的有：EVI、RVI、SAVI、GCVI、RDVI；平均方式优选的指数出现四次以上的是 GCVI、RDVI、NDVI、SAVI、EVI、MSR、NDre1、NDre2、NDVIre1。综合两种方式所选的指数，同样以出现五次以上的指数作为 ED 距离方法

优选指数集，包括 GCVI、EVI、RDVI、SAVI、MSR、NDre1、NDre2、NDVire1、gNDVI、RVI、NDVI、NDWI。所选水体指数只有 NDWI 一个。

4.1.3 SAD 距离

光谱角检测是比较两类数据相似程度的光谱对比方法，当 SAD 值越趋近于 1，两类的光谱就越相似。以 SAD 值小于 0.6 的指数出现的频次进行统计（附表 3），其指数出现频次统计如图 3-(c)所示，指数频次显示出较大差距，植被指数 MCARI、TCARI 不同于其他方法所选指数。进一步研究发现，TCARI 值是 MCARI 值的三倍，相比其他植被指数，对叶绿素浓度的变化十分敏感；红边指数 REP 出现最多，其次是 MTCI 指数。REP 指数可定量表征植物的叶绿素含量，而 MTCI 指数同样对叶绿素含量敏感。SAD 距离所选植被指数和红边指数均对植物叶片含的叶绿素有较好的表征。水体指数

仅选出 SWI，该指数能较好区分水体和阴影。

对 SAD 值分别进行加权和平均处理，优选结果如附表 6 所示。加权方式优选指数中出现四次以上的指数为：TCARI、MCARI、MCARI2、IRECI、EWI、GCVI、NDWI、MNDWI；平均方式优选指数出现次数四次以上的为：TCARI、MCARI、MCARI2、IRECI、CIre、RVI、GCVI、LSWI、EVI。将两种方式所选指数共同考虑得到 SAD 距离方法优选的指数集为 TCARI、MCARI、MCARI2、IRECI、CIre、RVI、GCVI、EWI、MNDWI、EVI、LSWI、NDWI。

4.1.4 RF-RFE 特征重要性排序

利用 RF-RFE 算法优选得到的前 12 位指数按重要性排序（表 5），不同的数理统计特征出现很多重复的指数，一定程度上说明这些指数在重要性排序中表现较好。

表 5 RF-RFE 算法特征重要性前 12 位特征指数集排序

Table 5 RF-RFE algorithm feature importance top 12 feature index set

数理统计特征	1	2	3	4	5	6	7	8	9	10	11	12
中值	RDVI	DVI	gNDVI	GI	NDWI	NDWI_B	MNDWI	RNDWI	EWI	CIre	VIgreen	TVI
平均值	VIgreen	NDVI	TVI	GCVI	GI	NDWI_B	MNDWI	RNDWI	EWI	CIre	RDVI	DVI
最大值	VIgreen	NDVI	RDVI	MSR	TVI	DVI	GI	MNDWI	UGWI	LSWI	RVI	NDWI_B
最小值	TCARI	RVI	DVI	SAVI	NDWI	NDWI_B	RNDWI	EWI	UGWI	SWI	EVI	CIre
标准差	VIgreen	TCARI	TVI	DVI	GI	RNDWI	EWI	LSWI	MCARI	MCARI2	MTCI	IRECI

（注：加粗数字越小，重要性程度越高）

其中，植被指数有 RDVI、DVI、GI、VIgreen、TVI，水体指数有 NDWI_B、MNDWI、RNDWI、

EWI，红边指数有 CIre。

4.1.5 Relief-F 特征重要性排序

依据 Relief-F 算法优选的指数按重要性排序（表 6），其中，中值、平均值和标准差特征除了指数重要性排序不同外，选择的指数均相同，并且部分指数也出现在最大值和最小值特征中。指数特征中表

现优异的指数分别为：水体指数中的 RNDWI、MNDWI、EWI，植被指数有 VIgreen、GI、MSR、GCVI、RVI，红边指数为 NDre1、NDVIre1、NDre2，反映出这些指数对于区分地类较为重要。另外还发现重要性排序第一位的总是水体指数。

表 6 Relief-F 算法特征重要性前 12 位特征指数集排序

Table 6 Relief-F algorithm feature importance top 12 feature index set

数理统计特征	1	2	3	4	5	6	7	8	9	10	11	12
中值	NDWI_B	VIgreen	GI	RNDWI	CIre	MNDWI	LSWI	NDVI	NDre2	NDre1	EWI	MSR
平均值	NDWI_B	VIgreen	GI	NDVI	CIre	RNDWI	MSR	MNDWI	LSWI	NDVIre1	NDre2	GCVI
最大值	MNDWI	LSWI	NDVI	EWI	VIgreen	gNDVI	MSR	NDWI	GI	UGWI	GCVI	NDre2
最小值	RNDWI	NDWI	SWI	RDVI	UGWI	SAVI	NDVIre1	TVI	MCARI2	EVI	RVI	MSR
标准差	RNDWI	MNDWI	VIgreen	GI	EWI	NDre1	MSR	NDVIre1	NDre2	GCVI	RVI	LSWI

（注：加粗数字越小，重要性程度越高）

4.2 数理统计特征的优选

样本的空间分布具有随机性特征，虽然在湿地分类中通常采用平均值的方式对样本进行处理，但哪一种统计特征更能代表样本的属性特征，目前为止没有相关的研究。为此本文研究分别根据平均值、中值、最大值、最小值和标准差进行各类特征的计算，在此基础上评价对属性特征优选的影响。

4.2.1 JM 距离

对表 5 基于 JM 距离优选的指数集分别采用随

机森林分类器进行分类，分类精度如附表 7 所示。

可以看出，中值和平均值特征总体精度和 Kappa 系数最高，均为 86.70%、0.840，标准差特征的最低。依据 JM 距离计算结果，指数数理统计特征的可分性能力可排序为平均值特征=中值特征>最小值特征>最大值特征>标准差特征。

4.2.2 ED 距离

对 ED 距离优选得到的指数集（附表 5）进行随机森林分类，分类精度结果如表 7。

表 7 ED 算法分类精度统计

Table 7 ED algorithm classification accuracy statistics

类别	B-a-2	B-b-2	B-c-1	B-d-1	B-e-1
----	-------	-------	-------	-------	-------

	PA/%	UA/%	PA/%	UA/%	PA/%	UA/%	PA/%	UA/%	PA/%	UA/%
草原	88.50	95.80	88.50	92.00	88.50	95.80	88.50	88.50	88.50	88.50
草甸	79.20	76.00	79.20	73.10	79.20	63.30	79.20	70.40	75.00	60.00
沼泽化草甸	76.00	70.40	72.00	72.00	64.00	69.60	72.00	75.00	56.00	73.70
沼泽	77.80	82.40	77.80	87.50	77.80	93.30	83.30	88.20	83.30	88.20
河流	100.00	100.00	100.00	100.00	95.80	88.50	100.00	100.00	95.80	100.00
裸地	100.00	100.00	100.00	96.30	92.30	96.00	96.20	100.00	92.30	85.70
总体精度/%	87.40		86.70		83.20		86.70		81.80	
Kappa 系数	0.849		0.840		0.798		0.840		0.781	

(注: 加粗数字为最大总体精度、Kappa 系数)

中值特征的总体精度与 Kappa 系数最高, 达到 87.40%、0.849; 标准差特征的最小; 平均值与最小值特征的相同。从 ED 方法看, 数理特征的可分性能力排序为: 中值特征>平均值特征=最小值特征>最大值特征>标准差特征。

4.2.3 SAD 距离

对 SAD 距离计算得到的指数集(附表 6)逐一进行随机森林分类, 分类精度结果附表 8。

将 SAD 距离方法优选的指数组合利用随机森林分类方法, 比较其精度后发现, 虽然所选出的指数与前两种方法选出的指数有较大不同, 但是平均值特征的总体精度依然能达到 87.40%, Kappa 系数达

0.848, SAD 距离结果显示不同统计特征可分性能力大小排序为: 中值特征>平均值特征>最大值特征=最小值特征>标准差特征。

4.2.4 RF-RFE 特征重要性排序

根据 RF-RFE 特征重要性排序选出指数集(表 5)的对应分类结果如附表 9 所示: 平均值特征的总体精度与 Kappa 系数最高, 达 87.40%、0.848, 而标准差特征的最低, 中值和最小值的相同。依据总体精度和 Kappa 系数的大小, 可将其排序为平均值特征>中值特征=最小值特征>最大值特征>标准差特征。

4.2.5 Relief-F 特征重要性排序

对 Relief-F 算法计算结果取前 12 位指数(表 6)进行随机森林分类得到结果见表 8。

表 8 Relief-F 算法分类精度统计

Table 8 Relief-F algorithm classification accuracy statistics

类别	E-a		E-b		E-c		E-d		E-e	
	PA/%	UA/%	PA/%	UA/%	PA/%	UA/%	PA/%	UA/%	PA/%	UA/%
草原	88.50	95.80	88.50	88.50	88.50	95.80	88.50	85.10	80.80	87.50

草甸	87.50	75.00	87.50	77.80	79.20	61.30	75.00	69.20	75.00	62.10
沼泽化草甸	68.00	73.90	68.00	73.90	64.00	64.00	72.00	75.00	64.00	72.70
沼泽	77.80	82.40	77.80	82.40	72.20	86.60	83.30	88.20	88.90	84.20
河流	100.00	100.00	100.00	100.00	95.80	95.80	100.00	100.00	100.00	100.00
裸地	100.00	96.30	100.00	100.00	92.30	100.00	96.20	100.00	96.20	100.00
总体精度/%	87.40		87.40		82.50		86.00		83.90	
Kappa 系数	0.849		0.849		0.790		0.832		0.807	

(注: 加粗数字为最大总体精度、Kappa 系数)

中值和平均值特征的总体精度和 Kappa 系数均为 87.40%、0.849, 最大值特征的最低, 比中值和平均值低了 4.90%、0.059。以数理统计特征的维度, 依据总体特征和 Kappa 系数的大小, 可将其排序为中值特征=平均值特征>最小值特征>标准差特征>最大值特征。

总体而言, 湿地信息提取最适合的数理统计特征为中值和平均值。

4.3 高寒湿地类型的区分能力

本文将沼泽、沼泽化草甸、草甸、河流、草原和裸地等土地覆被类型分别两两组合分析其可分性, 共有 15 组类型组合。

4.3.1 河流、裸地

河流与裸地因其截然不同的形状特征、下垫面环境, 是所有类别中分类精度最高的, 其制图精度和用户精度可达到 100%。由附表 1-3 可知, 能用于区分河流或裸地与其他类别的指数数目很多。

4.3.2 沼泽、沼泽化草甸、草甸、草原

由表 7-8 和附表 7-9 可知, 草原的制图精度和用户精度的中位数是 88.5%, 沼泽化草甸与草甸的制图与用户精度主要集中在 72%, 相比较低。沼泽相比于沼泽化草甸与草甸两类湿地过渡类型来说, 制图精度和用户精度均要高, 大部分情况集中在 83.3%。沼泽化草甸与沼泽、草原与草甸、草甸与沼泽化草甸这三种类型组合难以区分, 因此, 结合 JM 距离、ED 距离、SAD 距离计算结果, 尝试找出易于区分这三种组合类型的指数。具体做法如下: 依据 JM 距离、ED 距离、SAD 距离指数的中值特征计算结果, 选出沼泽化草甸与沼泽、草原与草甸、草甸与沼泽化草甸类型组合前 10% 的指数, 共计 25 个, 根据指数出现的频次排序, 选出前 10% 的指数为: MCARI2、NDWI、DVI、EVI、EWI、IRECI、MCARI、TCARI、UGWI。

5 讨论

5.1 分类指数集的构建

由于每种方法优选的指数不完全一致, 为确定对

于湿地分类最优的特征，我们分别取 Filter 方法（JM 距离、ED 距离、SAD 距离、Relief-F 算法）和 Wrapper 方法（RF-RFE 算法）计算结果中精度最高的指数集合，然后以这 5 个指数集合为基础，统计各个指数在 5 个集合中出现的次数，以众数为指标作为最终优选的指数，结果包括了 11 个指数。分别为植被指数 RDVI、NDVI、MSR、RVI、VIgreen，水体指数 RNDWI、NDWI、NDWI_B、MNDWI、EWI 和红边指数 CIre。

5.2 数理统计特征的选择

对优选出的 11 个指数，分别对 5 个统计特征进行随机森林分类，得到的总体精度和 Kappa 系数如表 9。

表 9 不同数理统计特征下分类精度统计

Table 9 Classification accuracy statistics under different mathematical statistics characteristics

数理统计特征	总体精度/%	Kappa 系数
中值	88.10	0.857
平均值	86.70	0.84
最大值	85.30	0.823
最小值	83.20	0.798
标准差	85.30	0.824

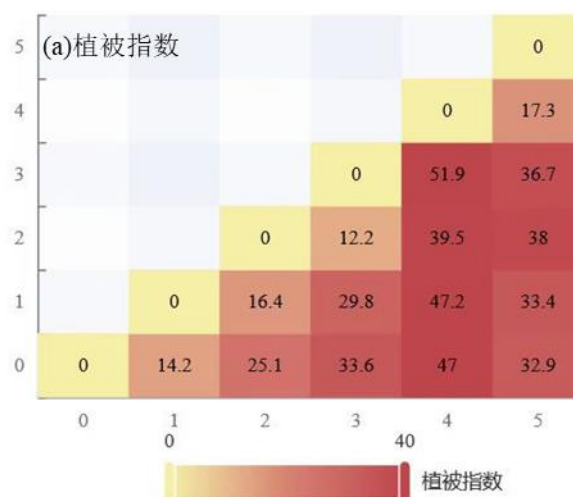
可以看出，中值特征的总体精度和 Kappa 系数均是最高，其次是平均值特征，和不同方法分别进行评价的结果一致（4.2 部分）。说明基于样本的中值或平均值可以代表样本的属性。

5.3 特征优选方法的比较

由于中值特征的精度以及分类的效果最好，因此用中值特征计算的精度及混淆矩阵来比较五种特优选的方法。可以发现 Filter 方法的 Relief-F 算法与 ED 距离算法的总体精度和 Kappa 系数均为 87.40%、0.849，JM 距离与 SAD 距离得到的总体精度和用户精度均为 86.70%、0.840，Wrapper 方法的 RFE-RF 算法的总体精度和用户精度为 86.00%、0.832，由此可知，Filter 方法的 Relief-F 与 ED 距离算法在本次高寒湿地分类指数特征研究中略胜一筹。

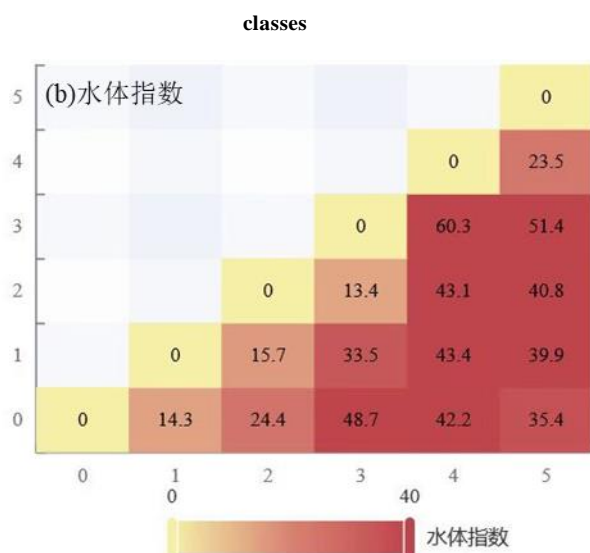
5.4 指数特征的重要性排序

为了能够更好的查看水体指数特征、植被指数特征与红边指数特征对湿地分类的贡献，利用 Filter 方法的 Z 检验方法对样本点的指数特征进行定量分析。依据 5.2 的结果，最终选择指数的中值特征进行计算。



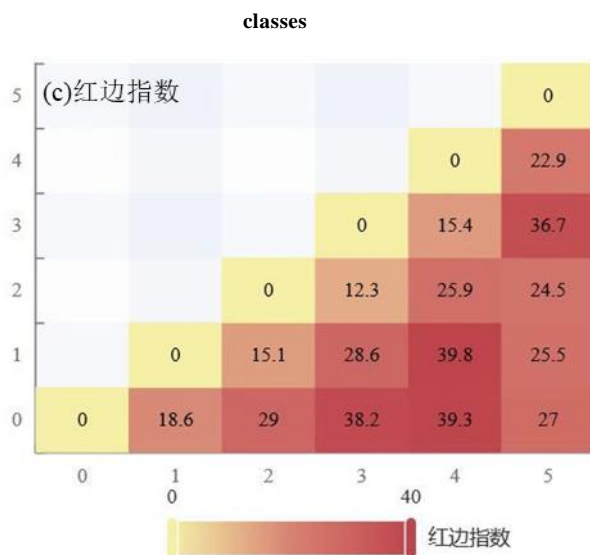
(a) 两两类之间植被指数显著性差异

(a) Significant difference in vegetation index between two



(b) 两两类之间水体指数显著性差异

(b) Significant difference of water body index between two



(c) 两两类之间红边指数显著性差异

(c) Significant difference in red edge index between two

classes

图 4 两两类之间的植被指数、水体指数和红边指数的显著性差异 (0-草原、1-草甸、2-沼泽化草甸、3-沼泽、4-河流、5-裸地)

Fig.4 Significant differences in vegetation index, water body index and red edge index between the two classes (0-grassland, 1-meadow, 2-swamped meadow, 3-swamp, 4-river, 5-bare land)

上图 4 表明了不同类型的指数特征对湿地的可分离度, Z 值越大, 两种湿地类型在此指数上的差

异度越明显。通过观察决定以 Z 值大于 37 为界限进行统计, 可得水体指数特征 Z 值最大为 60.3, 大于 37 的类型组合有 8 个; 植被指数特征 Z 值最大为 51.9, 大于 37 的类型组合有 5 个; 红边指数特征 Z 值最大为 39.8, 大于 37 的类型组合有 3 个。由此可知, 不同指数特征的特征重要性程度排序为: 水体指数特征>植被指数特征>红边指数特征。

5.5 分类特征的不确定性

沼泽化草甸与沼泽、草甸与沼泽化草甸、草原与草甸这三种类型组合的可分性较差, 而这三类组合的分类情况决定着湿地分类结果的好坏。从类型角度来看, 草原、草甸、沼泽化草甸、沼泽邻近两类型无论在植被长势还是下垫面水分含量都存在渐变的过程, 邻近两类型之间没有明显的地缘间隙。因此, 在用到的分类特征时是易混淆的。JM 距离、ED 距离、SAD 距离计算结果也证明了这一点, 后续需要深入研究上述类型的分类技术。

此外, SAD 距离方法所得指数集的随机森林分类精度与 JM 距离方法一样能达到 86.70%, 但是两种方法所选指数集却大相径庭, 接下来可以从指数的物理机理、波段组合方式方面进一步研究。

6 结论

本文基于 Sentinel-2 遥感影像, 以首曲高寒湿地为研究区, 通过 Filter 方法 (JM 距离、ED 距离、SAD 距离、Relief-F 算法)、Wrapper 方法

(RF-RFE 算法) 共同对植被指数特征、水体指数特征、红边指数特征进行分析, 并借助随机森林分类方法计算了所选指数集的精度和混淆矩阵, 最后利用了 Z 检验方法对三种指数特征进行了定量比较, 主要得到以下结论:

(1) 特征优选方法: 在所选的 Filter 和 Wrapper 方法中, Filter 方法的 ED 距离与 Relief-F 算法得出的指数集其分类精度高于其他方法, 精度最大相差 1.4%。说明 ED 距离与 Relief-F 算法在本次湿地分类研究中具有最好结果。其原因是 Wrapper 方法在特征选择中过分依赖聚类参数, 缺乏合适的评价准则评估不同特征子空间的特征子集, 所以采用 Filter 方法相对取得了较好的结果。

(2) 湿地分类优选指数: 植被指数 RDVI、NDVI、MSR、RVI、VGreen, 水体指数 RNDWI、NDWI、NDWI_B、MNDWI、EWI, 红边指数

CIre。具体来说, 水体指数在湿地分类中占重要地位, 尤其是 NDWI、MNDWI、NDWI_B, 植被指数在所选总指数中所占数量最多, 其中 NDVI 指数尤为重要, 红边指数虽没有前两类指数表现突出但是也不可或缺, 表现好的指数有 NDVIre1、NDre1、NDre2、CIre。通过 Z 检验对指数特征的定量分析, 可得出水体指数特征的重要性程度大于植被指数特征, 也大于红边指数特征; (3) 统计特征评价结果: 基于中值的总体精度最高, 达到 88.10%;

(4) 最易区分的是河流与裸地, 其次为沼泽与草原, 沼泽化草甸与草甸较为难分。沼泽化草甸与沼泽、草原与草甸、草甸与沼泽化草甸这三种类型组合较难区分, 可尝试如下指数进行分类:

MCARI2、NDWI、DVI、EVI、EWI、IRECI、MCARI、TCARI、UGWI。

参考文献(References)

- Broge N H, Leblanc E. 2001. Comparing prediction power and stability of broadband and hyperspectral vegetation indices for estimation of green leaf area index and canopy chlorophyll density. *Remote sensing of environment*, 76(2): 156-172[DOI: 10.1016/s0034-4257(00)00197-8]
- Breiman L. 2001. Random forest. *Machine Learning*, 45(1): 5-32 [DOI: 10.1023/A:1010933404324]
- Barnes EM, Clarke TR, Richards SR, Colaizzi, PD, Haberland J, Kostrzewski M, Waller P, Choi C, Riley E, Thompson T, Lascano RJ, Li H, Moran MS. 2000. Coincident detection of crop water stress, nitrogen status and canopy density using ground-based multispectral data. In: *Proceedings 5th International Conference Precision Agriculture*, Bloomington, USA, pp. 1-15.
- Birth G, McV. 1968. Measuring the color of growing turf with a reflectance Spectrophotometer. *Agron Journal*, 60, 640-643 [DOI: 10.2134/agronj1968.00021962006000060016x]
- Chen W J, Ding J L, Li Y H, et al. 2015. Extraction of water information based on China made GF-1 remote sense image. *Resources Science*, 37(6) : 1166-1172 (陈文倩, 丁建丽, 李艳华, 牛增懿. 2015. 基于国产 GF-1 遥感影像的水体提取方法. *资源科学*, 37(6) : 1166-1172) [CNKI:SUN:ZRZY.0.2015-06-009]
- CHANG W T, WANG H, NING X G, ZHANG H C. 2020.

- Extraction of Zhalong Wetlands Information based on Images of Sentinel-2 Red-edge Bands and Sentinel-1 Radar Bands. *Wetland Science*, 18(01): 10-19 (常文涛, 王浩, 宁晓刚, 张翰超. 2020. 融合 Sentinel-2 红边波段和 Sentinel-1 雷达波段影像的扎龙湿地信息提取. *湿地科学*, 18(01): 10-19) [DOI: 10.13248/j.cnki.wetlandsci.2020.01.002]
- Chen J M. 1996. Evaluation of vegetation indices and a modified simple ratio for boreal applications. *Canadian Journal of Remote Sensing*, 22(3): 229-242[DOI: 10.1080/07038992.1996.10855178]
- Duan J W, Zhong J S, Jiang L, Zhong M M. 2021. Extraction method of ultra-green water body index in post-rain flood area based on GF-2 image. *Geography and Geographic Information Science* (03) , 35-41(段纪维, 钟九生, 江丽, 钟淼淼. 2021. 基于 GF-2 影像的雨后洪涝区超绿水体指数提取方法. *地理与地理信息科学* , (03) : 35-41)[DOI: 10.3969/j.issn.1672-0504.2021.03.006]
- Daughtry C S T, Walthall C L, Kim M S, et al. 2000. Estimating corn leaf chlorophyll concentration from leaf and canopy reflectance. *Remote sensing of Environment*, 74(2): 229-239[DOI: 10.1016/s0034-4257(00)00113-9]
- Dash J, Curran P J. The MERIS terrestrial chlorophyll index. 2004. *International Journal of Remote Sensing*, 25 (23): 5403-5413 [DOI: 10.1080/0143116042000274015]
- Defries RS, Townshend JRG. 1994. NDVI-Derived Land Cover Classification at a Global Scale. *International Journal of Remote Sensing*, 15(17): 3567-3586[DOI: 10.1080/01431169408954345]
- Elavarasan, D., Vincent, P., Srinivasan, K., & Chang, C.Y. 2020. A Hybrid CFS Filter and RF-RFE Wrapper-Based Feature Extraction for Enhanced Agricultural Crop Yield Prediction Modeling. *Agriculture*, 10(9)[DOI: 10.3390/agriculture10090400]
- Frampton W J, Dash J, Watmough G, et al. 2013. Evaluating the capabilities of Sentinel-2 for quantitative estimation of biophysical variables in vegetation. *ISPRS journal of photogrammetry and remote sensing*, 82: 83-92 [DOI: 10.1016/j.isprsjprs.2013.04.007]
- F. A. Kruse, A. B. Lefkoff, and J. B. Dietz. 1993. Expert system-based mineral mapping in northern death valley, California/Nevada, using the airborne visible/infrared imaging spectrometer (AVIRIS). *Remote Sensing of Environment*, vol. 44, no. 2-3, pp. 309-336[DOI: 10.1016/0034-4257(93)90024-r]
- Gitelson A A, Kaufman Y J, Stark R, et al. 2002. Novel algorithms for remote estimation of vegetation fraction. *Remote sensing of Environment*, 80(1): 76-87[DOI: 10.1016/s0034-4257(01)00289-9]
- Gitelson AA, et al. 2003. Remote estimation of leaf area index and green leaf biomass in maize canopies. *Geophys Res Lett*, 30(5): 1248. [DOI:10.1029/2002GL016450]
- A. R. Huete. 1988. A soil-adjusted vegetation index (SAVI). *Remote Sensing of Environment*. 25:295-309[DOI: 10.1016/0034-4257(88)90106-X]
- Gitelson A A, Merzlyak M N. 1996. Signature analysis of leaf reflectance spectra: algorithm development for remote sensing of chlorophyll. *Journal of plant physiology*, 148(3-4): 494-500[DOI: 10.1016/s0176-1617(96)80284-7]
- Gitelson AA, Merzlyak M. Spectral reflectance changes associated with autumn senescence of *Aesculus hippocastanum* L. and *Acer platanoides* L. leaves. 1994. Spectral features and relation to chlorophyll estimation. *J. Plant Physiol*. 143, 286-292 [DOI: 10.1016/s0176-1617(11)81633-0]
- Gitelson A A, Keydan G P, Merzlyak M N. 2006. Three-band model for noninvasive estimation of chlorophyll, carot-enoids, and anthocyanin contents in higher plant leaves. *Geophysical research letters*, 33(11): 1-5 [DOI: 10.1029/2006gl026457]
- Guyot G, Baret F. 1988. Utilisation de la haute resolution spectrale pour suivre l'etat des couverts vegetaux. *Proceedings of the 4th international conference on spectral signatures of objects in remote sensing*. ESA SP-287, Assois, France: 1988: 279-286.
- Guyon, I, Elisseeff, A. 2003. An Introduction to Variable and Feature Selection. *Learn. Res*, 3, 1157-1182.
- G John, R. Kohavi, K. Pflieger. 1994. Irrelevant features and the subset selection problem. *The Eleventh Int'l Conf. Machine Learning*, 121-129.

- GAO B B. 2008. Research on the Status Quo of Wetland in Shouqu Nature Reserve of the Yellow River in Gansu Province and Its Protection Countermeasures, *Gansu Science and Technology*, 2008(12): 3-5(高斌斌. 甘肃黄河首曲自然保护区湿地现状及其保护对策研究. *甘肃科技*, 2008(12): 3-5)[DOI: 10.3969/j.issn.1000-0952.2008.12.002]
- HAO Y F, MAN W D, WANG J H, LIU M Y, ZHANG K. 2021. Wetland information extraction based on Relief-F algorithm and decision tree method. *Journal of Liaoning Technical University(Natural Science)*, 40(3): 225-233 (郝玉峰, 满卫东, 汪金花, 刘明月, 张阔. 2021. Relief-F 算法及决策树方法下的湿地信息提取. *辽宁工程技术大学学报(自然科学版)*, 40(3): 225-233) [DOI: 10.11956/j.issn.1008-0562.2021.03.006]
- HE J H, ZHANG T B, YI G H, BIE X J, LUO N, WANG Q, FAN W W. 2015. Wetland Information Extraction and Change Detection in Zoige Plateau Area Based on the EOS /MODIS. *Geomatics & Spatial Information Technology*, 38(09): 38-41 (何菊红, 张廷斌, 易桂花, 别小娟, 罗娜, 王强, 范微维. 2015. 基于 EOS/MODIS 若尔盖高原地区湿地信息提取及变化监测. *测绘与空间地理信息*. 38(09): 38-41) [DOI: CNKI:SUN:DBCH.0.2015-09-013]
- Huete, A., Didan, K., Miura, T., Rodriguez, E.P., Gao, X., Ferreira, L.G. 2002. Overview of the radiometric and biophysical performance of the MODIS vegetation indices. *Remote Sensing of Environment*, 83, 195-213[DOI: 10.1016/S00344257(02)00096-2]
- Haboudane D, Miller J R, Tremblay N, et al. 2002. Integrated narrow-band vegetation indices for prediction of crop chlorophyll content for application to precision agriculture. *Remote sensing of environment*, 81(2): 416-426[DOI: 10.1016/s0034-4257(02)00018-4]
- Korhonen L, Packalen P, Rautiainen M. 2017. Comparison of Sentinel-2 and Landsat 8 in the estimation of boreal forest canopy cover and leaf area index. *Remote Sensing of Environment*, 195: 259-274[DOI: 10.1016/j.rse.2017.03.021]
- Kira K. 1992. The feature selection problem: traditional methods and a new algorithm. *Proc. AAAI-92*.
- LIU J H, WANG F. 2021. A Semi-supervised Feature Selection Algorithm Based on Relief-F. *Journal of Zhengzhou University(Natural Science Edition)*, 53(01): 42-46+53(刘吉超, 王锋. 2021. 基于 Relief-F 的半监督特征选择算法. *郑州大学学报(理学版)*, 53(01): 42-46+53 [DOI: 10.13705/j.issn.1671-6841.2020196]
- L. Yu, H. Liu. 2003. Feature Selection for High-Dimensional Data: A Fast Correlation-Based Filter Solution. *Proc. 20th Int'l Conf. Machine Learning*, 856-863
- Masoud MahdianPari, Bahram Salehi, Fariba Mohammadimanesh, Saeid Homayouni, Eric Gill. 2019. The First Wetland Inventory Map of Newfoundland at a Spatial Resolution of 10 m Using Sentinel-1 and Sentinel-2 Data on the Google Earth Engine Cloud Computing Platform. *Remote Sensing*, 11(1) [DOI: 10.3390/rs11010043]
- McFeeters S K. 1996. The use of normalized difference water index(NDWI)in the delineation of open water features.*International Journal of Remote Sensing*, 17(7): 1425-1432[DOI: 10.1080/01431169608948714]
- M. Dabboor, S, Howell, M. Shokr, J. Yackel. 2014. The Jeffries-Matusita distance for the case of complex Wishart distribution as a separability criterion for fully polarimetric SAR data. *International Journal of Remote Sensing*, 35(19)[DOI: info:doi/10.1080/01431161.2014.960614]
- Moor H, Hylander K, Norberg J. 2015. Predicting climate change effects on wetland ecosystem services using species distribution modeling and plant functional traits. *Ambio*, 404(S1)[DOI: 10.1007/s13280-014-0593-9]
- M. Dash, H. Liu. Feature Selection for Classification. 1997. *Intelligent Data Analysis*, 1(3): 131-157 [DOI: 10.1016/S1088-467X(97)00008-5]
- M. Dash, K. Choi, P. Scheuermann, H. Liu. 2002. Feature Selection for Clustering-a Filter Solution.*Proc. Second Int'l Conf. Data Mining*, 115-122[DOI: 10.1109/ICDM.2002.1183893]
- Han, N, Wang, K, Yu, L, Zhang, XY. 2011. Integration of texture and landscape features into object-based classification for delineatingTorreyausing IKONOS imagery. *International Journal of Remote Sensing*, 2003-2033[DOI:

- 10.1080/01431161.2011.605084]
- Osmar A. Carvalho Júnior, Renato F. Guimarães, Alan R. Gillespie, Nilton C. Silva, Roberto A. T. Gomes. 2011. A New Approach to Change Vector Analysis Using Distance and Similarity Measures. *Remote Sensing*, 3(11)[DOI: 10.3390/rs3112473]
- Phan T N, Kuch V and Lehnert L W. 2020. Land Cover Classification using Google Earth Engine and Random Forest Classifier—The Role of Image Composition. *Remote Sensing*, 12: 2411 [DOI: 10.3390/rs12152411]
- Roujean J L, Breon F M. 1995. Estimating PAR absorbed by vegetation from bidirectional reflectance measurements. *Remote Sensing of Environment*, 51(3): 375-384[DOI: 10.1016/0034-4257(94)00114-3]
- Sun Y L, Zhang X, Shuai T, Shang K and Feng S N. 2015. Radiometric normalization of hyperspectral satellite images with spectral angle distance and Euclidean distance. *Journal of Remote Sensing*, 19(4) : 618-626(孙艳丽, 张霞, 帅通, 尚坤, 冯淑娜. 2015. 光谱角-欧氏距离的高光谱图像辐射归一化. *遥感学报*, 19(4): 618-626)[DOI: 10.11834/jrs.20154176]
- Shoko C and Mutanga O. 2017. Examining the strength of the newly launched Sentinel 2 MSI sensor in detecting and discriminating subtle differences between C3 and C4 grass species. *ISPRS Journal of Photogrammetry and Remote Sensing*, 129: 32-40 [DOI: 10.1016/j.isprsjprs.2017.04.016]
- Saeyes Y, Inza I, Larranaga P. 2007. A review of feature selection techniques in bioinformatics. *Bioinformatics*, 23(19): 2507-2517[DOI: 10.1093/bioinformatics/btm344]
- Tucker CJ. 1979. Red and photographic infrared linear combinations for monitoring vegetation. *Remote Sensing and Environment*, 8(2): 127-150[DOI: 10.1016/00344257(79)90013-0]
- White Lori, Ryerson Robert A. , Pasher Jon, Duffe Jason. 2020. State of Science Assessment of Remote Sensing of Great Lakes Coastal Wetlands: Responding to an Operational Requirement. *Remote Sensing*, 12(18) [DOI: 10.3390/rs12183024]
- Wu C, Niu Z, Tang Q, et al. 2008. Estimating chlorophyll content from hyperspectral vegetation indices: Modeling and validation. *Agricultural and forest meteorology*, 148(8): 1230-1241[DOI: 10.1016/j.agrformet.2008.03.005]
- Wu C W, Liang J H, Wang W, Li C S. 2017. Random Forest Algorithm Based on Recursive Feature Elimination. *Random Forest Algorithm Based on Recursive Feature Elimination*, 2017(21): 60-63(吴辰文, 梁靖涵, 王伟, 李长生. 2017. 基于递归特征消除方法的随机森林算法. *统计与决策*, 2017(21): 60-63)[DOI: 10.13546/j.cnki.tjyj.2017.21.014]
- Wang G X, Li Y S, Wang Y B, Chen L. 2007. Typical alpine wetland system changes on the Qinghai-Tibet Plateau in recent 40 years. *Acta Geographica Sinica*, (05): 481-491 (王根绪, 李元寿, 王一博, 陈玲. 近 40 年来青藏高原典型高寒湿地系统的动态变化. 2007. *地理学报*, (05):481-491) [DOI: 10.3321/j.issn:0375-5444.2007.05.004]
- XIE S Y, FU B L, LI Y, LIU Z L, ZUO P P, LAN F W, HE H C, FAN D L. 2021. Classification Method on Marsh Wetlands in Honghe National Nature Reserve based on Multi-dimensional Remote Sensing Images. *Wetland Science*, 19(01): 1-16(解淑毓, 付波霖, 李颖, 刘兆礼, 左萍萍, 蓝斐芜, 何宏昌, 范冬林. 2021. 基于多维度遥感影像的洪河国家级自然保护区沼泽湿地分类方法研究. *湿地科学*, 19(01): 1-16)[DOI: 10.13248/j.cnki.wetlandsci.2021.01.001]
- XUE P F, LI W L, ZHUG F, ZHOU H K, LIU C L, and YAN H P. 2021. Changes in the pattern of an alpine wetland landscape in Maqu County in the first meander of the Yellow River. *Chinese Journal of Plant Ecology*, 45(05): 467-475(薛鹏飞, 李文龙, 朱高峰, 周华坤, 刘陈立, 晏和飘. 黄河首曲玛曲县高寒湿地景观格局演变. *植物生态学报*, 2021, 45(05): 467-475)[DOI: 10.17521/cjpe.2020.0288]
- Xu X L, Liu J Y, Shao Q Q, Fan J W. 2008. The dynamic changes of ecosystem spatial pattern and structure in the Three-River Headwaters region in Qinghai Province during recent 30 years. *Geographical Research*, (04): 829-838+974 (徐新良, 刘纪远, 邵全琴, 樊江文. 2008. 30 年来青海三江源生态系统格局和空间结构动态变化. *地理研究*, (04): 829-838+974) [DOI:

- 10.11821/yj2008040011]
- Xu H Q. 2005. A study on information extraction of water body with the modified normalized difference water index (MNDWI). 2005. *Journal of Remote Sensing*, 9(5): 589-595(徐涵秋. 2005. 利用改进的归一化差异水体指数(MNDWI) 提取水体信息的研究. 2005. *遥感学报*, 9(5): 589-595) [DOI: 10.3321/j.issn:1007-4619.2005.05.012]
- Xiao, X. , Boles, S. , Liu, J., Zhuang, D., Froking, S., Li, C., Salas, W., Moore, B. 2005. Mapping paddy rice agriculture in southern China using multi-temporal MODIS images. *Remote Sensing of Environment*, 95, 480-492 [DOI: 10.1016/j.rse.2004.12.009]
- Yan P, Zhang Y J, Zhang Y. 2007. A study on information extraction of water system in semi- arid regions with the enhanced water index (EWI) and GIS based noise remove techniques. *Remote Sensing Information*, (6): 62-67. 闫霏, 张友静, 张元. 2007. 利用增强型水体指数(EWI) 和 GIS 去噪音技术提取半干旱地区水系信息的研究. *遥感信息*, (6): 62-67[DOI: 10.3969/j.issn.1000-3177.2007.06.015]
- Zheng Y, Wu B F and Zhang M. 2017. Estimating the above ground biomass of winter wheat using the Sentinel-2 data. *Journal of Remote Sensing*, 21(2): 318-328 (郑阳, 吴炳方, 张淼. 2017. Sentinel-2 数据的冬小麦地上干生物量估算及评价. *遥感学报*, 21(2): 318-328) [DOI: 10.11834/jrs.20176269]
- ZHANG L, GONG Z N, WANG Q W, JIN D D, WANG X. 2019. Wetland mapping of Yellow River Delta wetlands based on multi-feature optimization of Sentinel-2 images. *Journal of Remote Sensing*, 2019, 23(02): 313-326 (张磊, 宫兆宁, 王启为, 金点点, 汪星. Sentinel-2 影像多特征优选的黄河三角洲湿地信息提取. *遥感学报*, 2019, 23(02): 313-326) [DOI:10.11834/jrs.20198083]
- ZHANG S Q, ZHOU B G, SHI F F, et al. 2020. Study on Information Extraction Method of Alpine Wetland in Qinghai-Xizang Plateau based on Remote Sensing Data of GF-1 Satellite—Taking Maduo County for Example. *Plateau Meteorology*, 39(6): 1309-1317(张帅旗, 周秉荣, 史飞飞. 2020. 基于高分一号卫星遥感数据的青藏高原高寒湿地信息提取方法研究——以玛多县为例. *高原气象*, 39(6): 1309-1317) [DOI:10. 7522/j. issn. 1000-0534. 2019. 00131]

Research on Remote Sensing Feature Selection for Alpine Wetland Classification

Huo Xuanlin^{1,2}, Niu Zhenguo², Zhang Bo², Liu Linsong², Li Xia³

1. School of Earth Science and Resources, Chang'an University, Shaanxi 710054;

2. State Key Laboratory of Remote Sensing Science, Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing 100094, China;

3. School of Land Engineering, Chang'an University, Shaanxi 710054

Abstract: Alpine wetlands are one of the important surface cover types on the Qinghai-Tibet Plateau, and play a key role in water conservation, climate regulation, and biodiversity maintenance. Accurate and timely knowledge of the temporal and spatial distribution of alpine wetlands is necessary for wetland protection and management. The selection of remote sensing classification features plays a key

role in wetland mapping. Although spectral features, texture features, topographic features, etc. have all been involved in existing research, few studies have focused on spectral index features, and discussed their mathematical statistical features and feature selection methods in depth. Based on Sentinel-2 image data and taking the Shouqu Alpine Wetland Reserve as the research site, this paper aims to classify alpine wetlands from the aspects of mathematical statistical features, alpine wetland types, feature selection methods and selected feature sets combined with random forest classification algorithm. In-depth and comprehensive analysis and research on the spectral index characteristics of alpine wetlands, and a beneficial attempt to optimize the classification characteristics of alpine wetlands.

In this study, the Gansu Shouqu Alpine Wetland Reserve was used as the research area, and the classification characteristics (spectrum, vegetation index, red edge index and water body index) were obtained based on Sentinel-2 data. Using Filter and Wrapper feature selection methods including Jeffries-Matusita distance, spectral angular distance (SAD), Euclidean Distance (ED), RF-RFE algorithm and Relief-F algorithm are used to optimize the above features, and use the Z test method for quantitative evaluation.

The research shows that: (1) Among all the categories of alpine wetlands involved in the classification, rivers and bare land are the easiest to distinguish, followed by grasslands and swamps, and swampy meadows and meadows are more difficult to distinguish. MCARI2, NDWI, DVI, EVI, EWI, IRECI, MCARI, TCARI, and UGWI indices can be used to distinguish between swamps, swampy meadows, meadows, and grasslands adjacent to each other; (2) Contribution of different index characteristics to wetland information extraction In terms of degree, water body index characteristics > vegetation index characteristics > red edge index characteristics; (3) From the perspective of feature optimization methods, the ED distance algorithm and the Relief-F algorithm in the Filter method are outstanding; (4) Finally, the most suitable alpine wetland information extraction is selected. The indices are RDVI, NDVI, MSR, RVI, VIgreen, RNDWI, NDWI, NDWI_B, MNDWI, EWI, CIre; (5) From the mathematical statistics of different classification features, the median feature has the best classification result, followed by the average value feature.

By means of multi-dimensional analysis, we give detailed results from feature optimization methods, wetland classification optimization index, statistical feature evaluation, and categories involved in classification in alpine wetlands. To the best of our knowledge, this study provides a transferable and universal method and idea for the selection of characteristic variables for wetland information extraction.

Key words: wetland classification; alpine wetland; feature selection; Qinghai-Tibet Plateau; Sentinel

Supported by The Second Research Project of Comprehensive Scientific Expedition to the Qinghai-Tibet Plateau (2019QZKK0106); Research on Chinese Wetland Classification Method Based on Time Series Remote Sensing Data (41971390).