

# 基于稀疏网络的可见光/近红外反射光谱 土壤有机质含量估算

冉思<sup>1,2</sup>, 丁建丽<sup>1,2\*</sup>, 葛翔宇<sup>1,2</sup>, 刘博华<sup>1,2</sup>, 张钧泳<sup>1,2</sup>

<sup>1</sup> 新疆大学资源与环境科学学院, 新疆 乌鲁木齐 830046;

<sup>2</sup> 新疆大学绿洲生态教育部重点实验室, 新疆 乌鲁木齐 830046

**摘要** 采集艾比湖湿地 89 个典型样点和土壤实测光谱数据, 对所测土壤光谱进行一阶微分变换预处理, 采用连续投影算法(SPA)、主成分分析(PCA)和稀疏自编码(SAE)对光谱数据进行特征提取, 结合偏最小二乘回归与 BP (Back Propagation)神经网络构建 SOM 估算模型。实验结果表明, SAE 方法能够有效对光谱进行压缩; 相比于 PLSR 模型, BP 模型能够较好地处理光谱中复杂的非线性信息; SAE-BP 方法在估算 SOM 中取得的精度最高。网络模型的建模方式能够显著提高 VIS-NIR 光谱反演土壤有机质模型的稳定性和精度, 当面对光谱中复杂的非线性问题时, 具有很强的解析力和较好的模型稳健性, 为使用 VIS-NIR 数据估算 SOM 提供一种新思路。

**关键词** 遥感; 土壤有机质; 可见-近红外光谱; 稀疏自编码; BP 神经网络

中图分类号 O433 文献标志码 A

doi: 10.3788/LOP57.242803

## Estimation Method of VIS-NIR Spectroscopy for Soil Organic Matter Based on Sparse Networks

Ran Si<sup>1,2</sup>, Ding Jianli<sup>1,2\*</sup>, Ge Xiangyu<sup>1,2</sup>, Liu Bohua<sup>1,2</sup>, Zhang Junyong<sup>1,2</sup>

<sup>1</sup> College of Resources & Environmental Science, Xinjiang University, Urumqi, Xinjiang 830046, China;

<sup>2</sup> Key Laboratory of Oasis Ecology, Xinjiang University, Urumqi, Xinjiang 830046, China

**Abstract** This research presents a novel approach for using VIS-NIR spectroscopy for soil organic matter (SOM) estimation. Soil spectrum data is collected from 89 samples retrieved from the Aibi Lake wetland. The samples are measured using a first-order differential transformation achieved through a continuous projection algorithm, a principal component analysis, and a sparse auto-encoder (SAE). The extracted data is then combined with a partial least squares regression (PLSR) and backpropagation (BP) neural network for the purpose of building a SOM estimation model. Experimental results show that the SAE method is able to effectively compress the spectrum. The BP model is shown to handle the complex and nonlinear information of the spectrum better than the PLSR model. Meanwhile, the SAE-BP method has the highest accuracy for estimating SOM. The network model is shown to significantly improve the stability and accuracy of the vis-NIR spectrum inversion of the SOM model. This model shows a robust and strong analytical power when faced with complex nonlinear problems in the spectrum.

**Key words** remote sensing; soil organic matter; visible-near infrared spectroscopy; sparse self-encoding; BP neural network

**OCIS codes** 280.4750; 140.4780; 140.3390; 140.3945

## 1 引言

土壤有机质(SOM)作为土壤的重要养分来源

之一,能够提高土壤肥力、促进植物的生长以及实现农业的可持续发展,而且及时有效地掌握 SOM 含量的变化,对湿地的保护与维护有重要的指导意

收稿日期: 2020-02-07; 修回日期: 2020-02-28; 录用日期: 2020-03-06

基金项目: 国家自然科学基金(41771470, 41961059)、新疆教育厅自然科学基金重点项目(XJEDU2018IOO8)

\* E-mail: watarid@xju.edu.cn

义<sup>[1-3]</sup>。传统的化学测定方法有较高的测量精度,但需要消耗大量的时间成本且价格昂贵,在实际应用中难以大规模使用。由于 SOM 在可见光和近红外(VIS-NIR)光谱中有着独特的吸收特性,因此常用来估算 SOM 含量。VIS-NIR 光谱技术作为研究已超过二十年的成熟技术,在估算 SOM 含量的准确性和经济效益上已得到科研工作者的广泛认可<sup>[4-6]</sup>。土壤光谱与 SOM 含量之间通常呈非线性相关,而光谱中包含若干噪声,从中探明与 SOM 有关的敏感变量仍存在一定难度,因此对土壤光谱进行特征变量的提取和干扰信息的剔除是模型具有较高准确性的重要保证<sup>[7-9]</sup>。

在 VIS-NIR 光谱数据的研究过程中,光谱的预处理方式与模型的构建是提高预测精度的关键。Hong 等<sup>[10]</sup>采用了不同的分数阶导数对土壤反射率光谱进行预处理,发现随着导数阶数的增加,基线漂移的现象和重叠峰逐渐消失,但更容易受到噪声的干扰。章海亮等<sup>[11]</sup>应用了遗传算法与连续投影算法对波段进行优化,结合偏最小二乘回归建立有机质预测模型,结果表明对原始波段进行筛选优化能够显著提高模型的预测精度。栾福明等<sup>[12-13]</sup>通过相关分析和主成分分析(PCA)选择了特征变量,结合 BP(Back Propagation)神经网络建立更为简洁的模型,从而提高了模型的预测性能。国内外学者采用不同的变量筛选与建模方法对 SOM 含量进行了大量的研究并取得了显著的成果,但这些模型难以充分探讨土壤光谱之间的线性与非线性关系,无法提取包含深层特征的光谱。由于模型收敛具有不稳定,常存在诸如局部最优和过拟合等问题<sup>[14]</sup>,如何准确提取 VIS-NIR 光谱中的信息显得尤为重要。

深度学习算法在机器学习领域中带来了一系列突破,该算法可以自动学习和提取光谱数据中的固有和深层次特征<sup>[15]</sup>。因采用深度学习算法估算 SOM 的研究较少,为此提出一种新颖的深度学习算法以反演 SOM。深度学习的核心是神经网络,如深度神经网络、卷积神经网络和自动编码器神经网络。与有监督的深度神经网络不同,自动编码器是一种无监督的特征学习神经网络,仅使用几层网络就可以从数据中提取特征<sup>[16]</sup>,已成功应用于医学预测<sup>[17]</sup>和图像处理<sup>[18]</sup>等领域。基于此,本文针对 SOM 提出一种新的稀疏自编码-BP(SAE-BP)网络模型,以期能够最大程度地提高模型的预测性能及精度,为今后研制土壤光谱传感器提供理论依据。

## 2 数据与方法

### 2.1 土壤样品的采集

研究区域位于新疆北部的艾比湖湿地,地处亚欧大陆腹地(43° 38' N—45° 52' N, 79° 53' E—85° 02' E),该区域的土壤质地以砂土、壤土、黏壤土和黏土为主。实验共采集 103 个样品,土样均匀分布于艾比湖湿地的周围,将样品带回实验室之后自然风干,研磨后过 0.5 mm 孔径的筛子过筛,采用重铬酸钾-硫酸溶液加热法来测定 SOM。

### 2.2 光谱数据的采集

采用美国 ASD 公司的 ASD Field Spec®3 HR 便携式光谱仪(波长范围为 350~2500 nm)来测定土壤的光谱反射率数据,重采样间隔为 1 nm,每个土壤数据累计得到 2151 条波段。光谱测定过程是在暗室中进行的,光源为 50 W 功率的卤化灯,探头的入射角度为 15°,探头距离样品表面为 10 cm,光源距离样品表面为 50 cm,每次反射率的测定均使用白板进行标定,每个样品均测量 5 次,取其平均值作为该样品的光谱反射率。实际的测定过程中,不同传感器之间的响应精度不同,使得位于两端的光谱数据易混入噪声信息,每份土样均对 350~400 nm 和 2401~2500 nm 波段的光谱进行剔除。

### 2.3 光谱数据的预处理及异常数据的剔除

为了突出原始数据的光谱信息、提高信噪比以及消除高频随机噪声对模型的影响,采用 SG(Savitzky-Golay)平滑以及微分处理对原始光谱数据进行预处理,其中微分处理包含一阶微分处理(FDR)和二阶微分处理(SDR)。SG 平滑可以提高光谱的平滑性,降低噪声的干扰;光谱微分处理能够部分消除外部环境的干扰,提高灵敏度与光谱的分辨率。此外,为了减少异常样本对模型性能产生干扰,使用 PCA 与马氏距离的结合方法对异常样本进行剔除<sup>[19]</sup>。

### 2.4 变量筛选方法

采用多种光谱变量筛选方法对原始光谱进行筛选,并结合 PLSR(Partial least squares regression)和 BP 网络模型来预测 SOM,目标算法主要包括连续投影算法(SPA)、PCA 与 SAE。

SPA 是一种使矢量空间共线性最小化的前向变量选择算法,其在向量空间中执行简单的投影操作以获得共线性最小的有用变量子集,能够在有效消除变量间共线性的同时得到最低限度的冗余信息的变量组合,从而实现在较低的模型复杂度下以最

大限度地获取解释信息。SPA 的变量选择原则为新选择的变量是在原选择变量的正交子空间上选择具有最大投影值的变量,其可以在方均根误差(RMSE)最小的基础上确定最优的初始变量和变量数量,因此常用于光谱特征波长的筛选。PCA 是一种常用的数据压缩算法,通过变换可以得到具有相同变量数量的新变量。这些新变量是原变量的线性组合且彼此正交,包含的信息不重叠,进而消除变量之间的多重共线性。理论上,采用 PCA 可以获得主成分的维度与原始变量数据相同,但由于前几个主成分的贡献率较大,则只需保留几个贡献较大的主成分,即可保留原始数据中大部分的信息。

自编码器(AE)尝试学习一个  $H_{w,b}(x) \approx x$  的函数,而实际输出变量  $\hat{x}$  近似于输入变量  $x$ ,其中  $w$  为权重,  $b$  为偏置。在 AE 的网络结构中,输入的神经元个数与输出个数相同,输入的数据与输出数据近似相等,并采用反向传播算法对其进行训练。当隐藏层的神经元数量少于输入数量时,编码器可以学习到数据中隐含的特征,从而达到数据压缩的效果。AE 的网络结构一般由三层组成,即输入层、隐藏层和输出层,结构如图 1 所示,其中  $n = \{1, 2, \dots, 200\}$ ,  $a$  为隐含层需要压缩的维度,若将一个样本压缩到 10 个特征,则  $m$  值为 10。

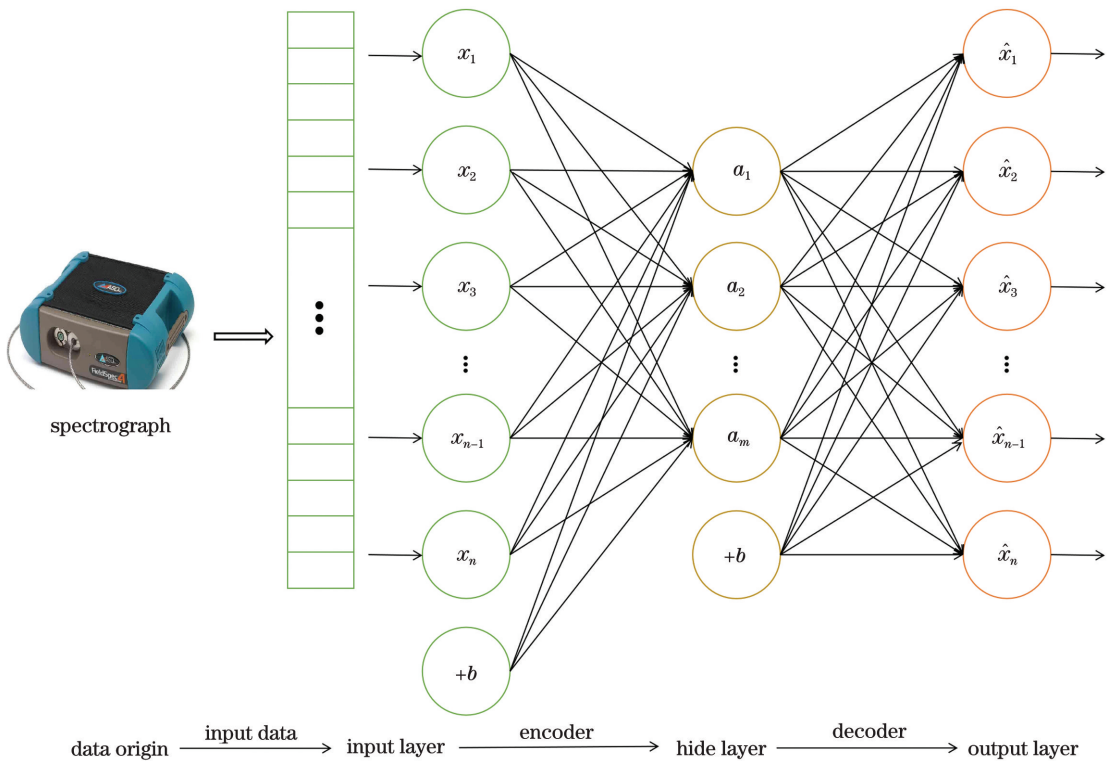


图 1 AE 网络结构

Fig. 1 AE network structure

SAE 是 AE 的改进,即在 AE 的优化目标函数中添加额外的惩罚因子,用来限制隐藏层中被激活的节点数,通过寻找一组超完备基来高效地表示数据的一种无监督学习方法,表达式为

$$L = \sum_{j=1}^S \rho \ln(\rho / \bar{\rho}_j) + (1 - \rho) \ln[(1 - \rho) / (1 - \bar{\rho}_j)], \quad (1)$$

式中:  $\rho$  为稀疏性参数;  $S$  为隐藏层神经元的总数;  $\bar{\rho}_j = \sum_{i=1}^m a_j^{(2)} x_i / m$ ,  $i$  为输入层中的神经元个数,  $j$  为隐藏层中的神经元个数,  $m$  为输入层神经元的总数,  $a_i$  为在给定输入层第  $i$  个维度的情况下,自编码

神经网络隐藏神经元第  $j$  个维度的激活度,上标(2)为第二层的隐藏层。总的代价函数可表示为

$$L = \min \sum_{i=1}^P (x_i - \hat{x})^2 + \beta \sum_{j=1}^S \{ \rho \ln(\rho / \bar{\rho}_j) + (1 - \rho) \ln[(1 - \rho) / (1 - \bar{\rho}_j)] \}, \quad (2)$$

式中:  $P$  为输入层神经元的总数;  $\beta$  为惩罚因子的权重。

### 2.5 模型评价

实验过程中,分别建立线性(多元线性回归模型和贝叶斯线性回归)和非线性(反向传播神经网络模型)两种不同的模型来预测 SOM。模型评价指标包

括决定系数( $R^2$ )、RMSE、剩余预测偏差(RPD)以及四分位数间隔(RPIQ)。 $R^2$ 用来表示模型的拟合程度, $R^2$ 值越接近于1,模型的拟合效果越好;RMSE用来表示模型的估算能力, RMSE值越小越好;RPD用来表示模型预测的准确性,通常来说当 $x_{RPD} < 1.4$ 时,意味着模型的量化能力较弱,当 $1.4 \leq x_{RPD} \leq 2.0$ 时,表示模型的量化能力尚可,当 $x_{RPD} > 2.0$ 时,表示模型的量化能力极强;RPIQ为四分位间距与RMSE的比<sup>[20]</sup>。采用 Kennard-Stone 算法来划分数据集,选取 59 个样本点作为训练集,30 个样本点作为验证集,分别建立预测模型。

### 3 结果与分析

#### 3.1 SOM 统计信息

实验过程中不可避免地存在误差,采用 PCA 与马氏距离的结合方法对异常样品进行剔除,共得到 89 个样本可以用于分析建模,如图 2 所示。

表 1 为 SOM 数据集的统计信息。从表 1 可以看到,SOM 含量在  $0.40 \sim 45.47 \text{ g} \cdot \text{kg}^{-1}$  之间,训练数据集的标准差(SD)和变异系数(CV)分别为

表 1 SOM 的统计信息

Table 1 SOM statistics

Dataset	Count	Min	Middle	Max	Mean	SD	CV
Whole dataset	89	0.40	12.56	45.47	13.07	8.08	0.62
Calibration dataset	59	0.40	12.56	45.47	12.67	6.98	0.55
Validation dataset	30	0.56	12.72	42.54	13.84	9.99	0.76

#### 3.2 光谱特征分析

图 3 为 89 个 SOM 样本的 VIS-NIR 光谱反射率曲线,这对于定性和定量研究 SOM 含量具有重要意义。从图 3(a)可以看到,原始光谱由三个吸收

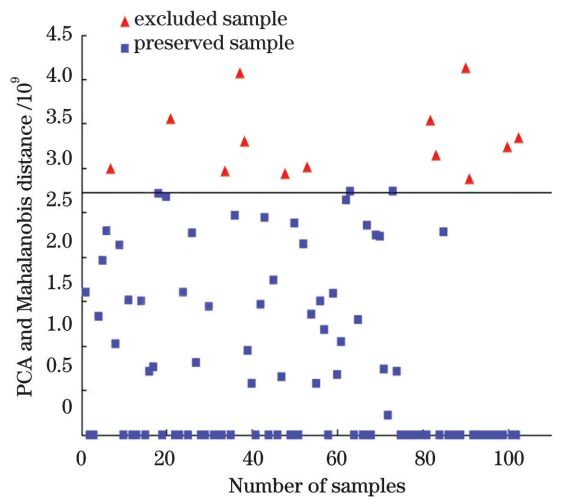


图 2 PCA 与马氏距离的结合方法剔除异常样本的结果  
Fig. 2 Results of combined method of PCA and Mahalanobis distance to eliminate abnormal samples

6.98 和 0.55,验证数据集的 SD 和 CV 分别为 9.99 和 0.76;高变异性即  $0.50 < x_{CV} < 1.00$ ,表明 SOM 含量在研究区域中是可变的,较大的土壤变异有助于提高校准模型的预测能力;训练集 SOM 含量范围覆盖验证集可以保证模型的合理估计。鉴于此结果,该模型可以提供良好的预测性能。

峰组成,分别在 1450,1950,2200 nm 附近,并且在可见光范围(400~760 nm)中,反射光谱呈现递增的形式,在近红外范围(780~2400 nm)中,由于光谱吸收带宽且重叠,这种趋势不太明显。从图 3(b)

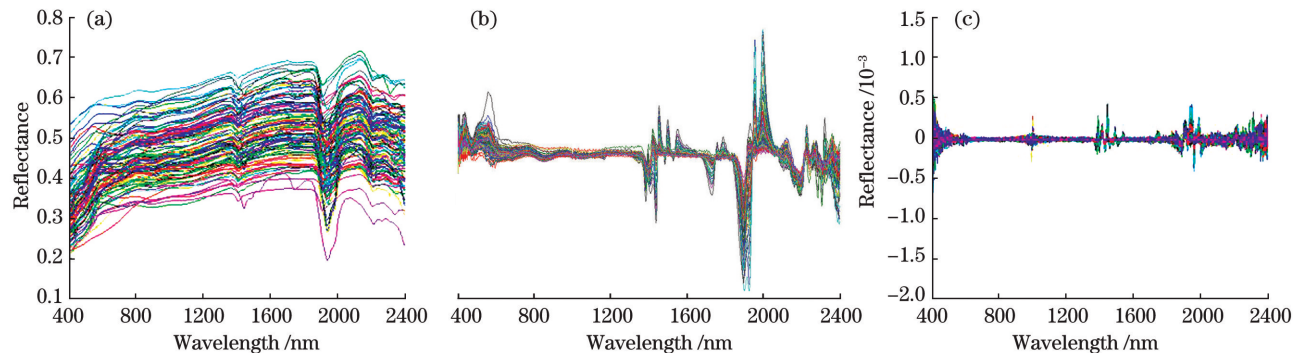


图 3 SOM 与不同 VIS-NIR 光谱的反射率曲线。(a)原始光谱;(b)一阶光谱;(c)二阶光谱

Fig. 3 Reflectance curves of SOM and different VIS-NIR spectra. (a) Original spectrum;

(b) first-order spectrum; (c) second-order spectrum

可以看到,在 1450,1950,2200 nm 附近的吸收峰,其反射能力相对增强。从图 3(c)可以看到,随着导数阶数的增加,多数反射率值逐渐接近于零,这表明重叠峰和基线漂移已被消除,此外在 VIS-NIR 波段出现许多小峰,这表明高阶导数的光谱更容易受到光谱噪声的干扰。

SOM 与不同 VIS-NIR 光谱的相关系数如图 4 所示。原始光谱与 SOM 之间的相关性较弱,在整个 VIS-NIR 波段,相关系数的波动范围为  $-0.15 \sim 0.20$ ,除了在  $400 \sim 1150$  nm 的波长范围内波动较大外,总体上曲线趋于平缓,如图 4(a)所示。当波长范围为  $400 \sim 850$  nm 和  $1650 \sim$

$2150$  nm 时,相关性波动较强,当波长为  $826$  nm 时,相关性达到最大,即  $r_{\max} = 0.64$ ,如图 4(b)所示。受到光谱噪声干扰等原因的影响,相关性在  $\pm 0.3$  的范围内反复波动,如图 4(c)所示。总体上,VIS 区域中 SOM 反射光谱之间的相关系数值大于 NIR 区域。此外,分数阶算法可以明显提高某些特定波长下 SOM 与分数阶光谱之间的相关性<sup>[21]</sup>。经过对比发现,FDR 的预处理效果最好,即  $r_{\max} = 0.64$ ,SDR 的预处理效果较差,故后续研究仅保留 FDR 的结果,应选 475 个波段 ( $r > 0.25$ ,显著性检验阈值  $Y = 0.01$ ) 作为研究的子集。

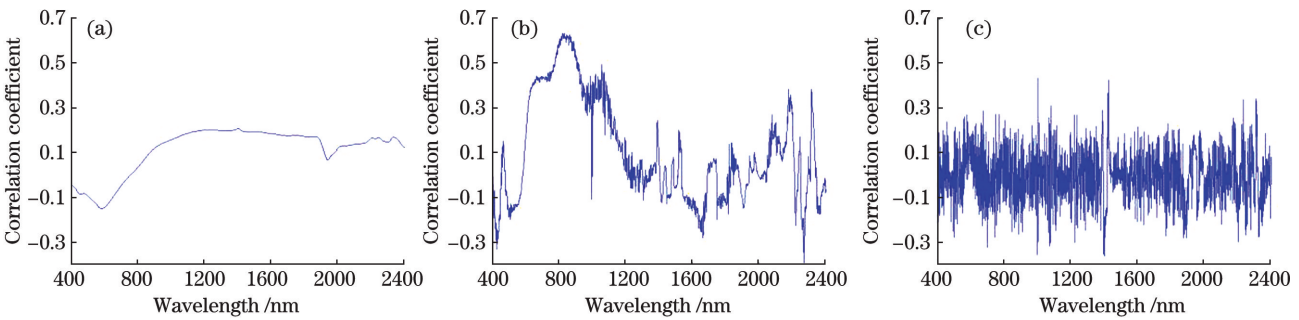


图 4 不同 VIS-NIR 光谱与 SOM 的相关性曲线。(a)原始光谱;(b)一阶光谱;(c)二阶光谱  
Fig. 4 Correlation curves of different VIS-NIR spectra and SOM. (a) Original spectrum;  
(b) first-order spectrum; (c) second-order spectrum

### 3.3 光谱特征变量的提取

为了获得具有稳健预测能力和少量输入变量的最优模型,采用三种变量选择方法(SPA、PCA 和 SAE)来筛选光谱特征,从而估算 SOM。

#### 3.3.1 PCA 和 SPA 方法提取特征变量

利用 SPSS 22.0 软件对一阶导数光谱进行 PCA 处理,选取特征值大于 1 的因子,结果如图 5 所示。从图 5 可以看到,特征值大于 1 的共有 18 个主成分分量,其中前 10 个主成分分量的累计贡献率

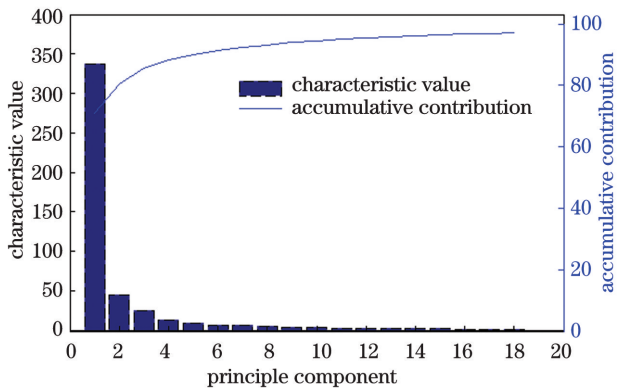


图 5 PCA 的特征值和贡献率

Fig. 5 PCA eigenvalues and contribution rate

为 94.541%,一般选取累计贡献率范围为 85%~95%的特征值所对应的主成分,因此选取前 10 个主成分来反映一阶导数光谱的信息。

图 6 为使用 SPA 变量筛选算法得到的特征波段 RMSE 曲线,即光谱变量筛选结果。从图 6 可以看到,随着筛选变量数量的增加,RMSE 值呈波动下降的趋势,当选取变量达到 84 个时,RMSE 达到

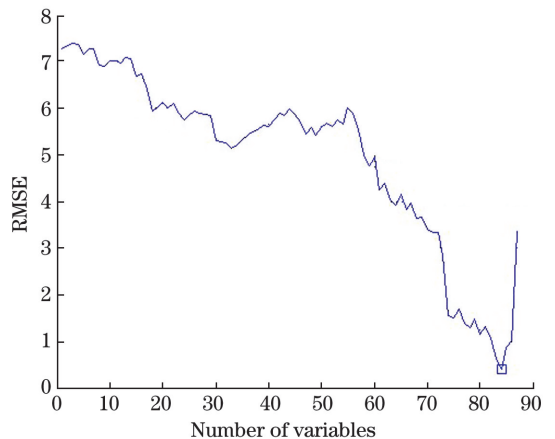


图 6 特征波段的 RMSE 曲线

Fig. 6 RMSE curve of characteristic band

最低值,即 0.398,而 SPA 算法筛选的光谱变量占一阶导数光谱变量的 18%。

### 3.3.2 SAE 方法提取特征变量

深度学习的学习能力受到参数的影响,准确选择参数对于模型的构建来说十分重要。传统的神经网络模型中,通常采用 Sigmoid 和 Tanh 函数作为激活函数,但是对于网络层数较多的神经网络往往存在计算量大、反向传播过程中梯度消失及信息丢失等问题,因此选择 ReLU 作为激活函数,可以减少运算量且降低过拟合发生的概率,优化器选用 Adam,损失函数采用均方对数误差函数。一次训练所选取的样本数设置为 30,迭代次数为 100 次,学习率为 0.001。选择 475 个一阶导数光谱作为输入层数据,将数据压缩至 200 个,再转换到 40 维。为了选择最优的维度,在网络的第 4 层将数据维度分别压缩至 4,6,8,10,12,14,16,18 维进行实验,最后特征变量通过解码层依次重建为 475 个光谱波段,SAE 网络结构及参数如图 7 所示。

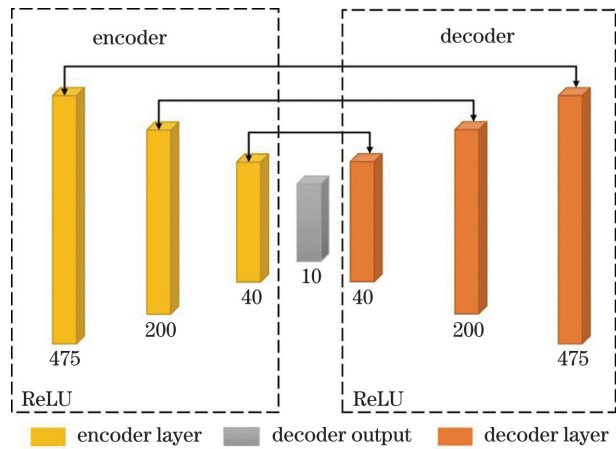


图 7 SAE 网络结构及参数

Fig. 7 SAE network structure and parameters

### 3.3.3 变量筛选方法的对比分析

变量筛选方法与 SOM 中  $R^2$  的关系,如图 8 所示,其中 IQR 为四分位差。从图 8 可以看到,在三种变量筛选方法中,SPA 方法所选择的变量数高达 84 个,导致变量与 SOM 中  $R^2$  之间的波动较大,PCA 方法和 SAE 方法中各变量之间的  $R^2$  值较为接近。相比于 SPA 方法和 PCA 方法,SAE 方法的  $R^2$  值最高,最大  $R^2$  值从大到小的顺序为 SAE、SPA、PCA,最小  $R^2$  值从大到小的顺序为 SAE、PCA、SPA, $R^2$  的中间值从大到小的顺序为 SAE、PCA、SPA,SAE 方法对于光谱的降维效果明显优于 SPA 方法和 PCA 方法。

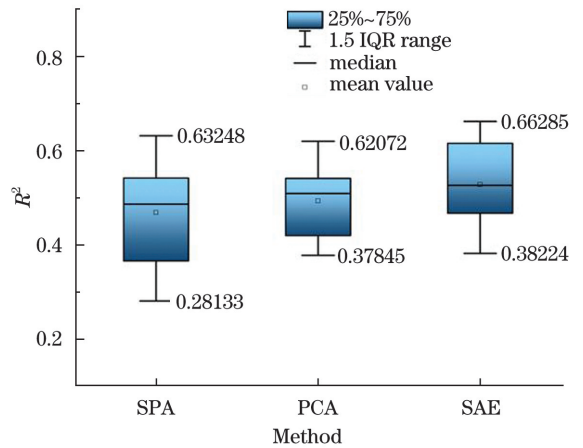


图 8 变量筛选方法与 SOM 中  $R^2$  的关系

Fig. 8 Relationship between variable selection method and  $R^2$  in SOM

### 3.4 模型的构建及对比分析

选择有效的变量能够大幅度提升模型的预测能力,采用 SAE、PCA 和 SPA 方法挑选不同的特征变量,并结合 PLSR 和 BP 神经网络来建立预测模型。表 2 为模型建模集与验证集的统计信息。从表 2 可以看到,SAE-BP 模型在建模集 ( $R^2 = 0.85, x_{RMSE} = 3.38$ ) 和验证集 ( $R^2 = 0.82, x_{RMSE} = 3.53, x_{RPD} = 2.32, x_{RPIQ} = 2.78$ ) 上的预测准确性最好;PCA-BP 模型在测试集上的预测精度 ( $R^2 = 0.76, x_{RMSE} = 4.08, x_{RPD} = 2.01, x_{RPIQ} = 2.40$ ) 略次于 SAE-BP 模型。对于 59 个训练样本与 30 个验证样本的结果预测,SPA-PLSR 模型在测试集上表现出最差的估算能力 ( $R^2 = 0.53, x_{RMSE} = 5.56, x_{RPD} = 1.47, x_{RPIQ} = 1.70$ )。相比于 PCA-PLSR 模型,SPA-PLSR 模型在测试集上的  $R^2$  降低了 13.11%,RMSE 增加了 9.45%,SAE-PLSR 模型在测试集上的  $R^2$  与之持平,RMSE 则增加了 2.95%;相比于 SAE-BP 模型,SPA-BP 和 PCA-BP 的  $R^2$  分别降低了 13.42% 和 7.32%,RMSE 则增加了 25.50% 和 15.58%。实验

表 2 不同特征变量筛选与建模方法的模型精度

Table 2 Model accuracy of different feature variable selections and modeling methods

Model	Extract variable	Calibration		Prediction			
		$R^2$	RMSE	$R^2$	RMSE	RPD	RPIQ
PLSR	SPA	0.64	5.02	0.53	5.56	1.47	1.70
	PCA	0.67	4.81	0.61	5.08	1.61	2.09
	SAE	0.66	4.90	0.61	5.23	1.58	2.07
BP	SPA	0.75	4.36	0.71	4.43	1.85	2.21
	PCA	0.80	3.90	0.76	4.08	2.01	2.40
	SAE	0.85	3.38	0.82	3.53	2.32	2.78

结果表明,基于 BP 模型的三种变量筛选方法的预测精度均高于 PLSR 模型,说明 SAE 方法能够有效地对 VIS-NIR 光谱进行压缩提取,且 SAE 方法与 BP 神经网络相结合的模式可以较好地估算 SOM。

## 4 讨 论

光谱变量的筛选是研究土壤 VIS-NIR 光谱的重要步骤,其可以有效去除光谱中的冗余信息,提高模型的预测精度。SPA 方法在 PLSR 模型和 BP 模型中的效果都较差,但其在很大程度上可以避免光谱信息的重叠,当筛选变量时倾向于选择不稳定的波段并剔除一些重要的相关波段,使得部分波段的信息缺失,从而影响模型的预测精度<sup>[22-23]</sup>。此外,较长的计算时间和所选波长的数量(不能大于校准样品的数量)是 SPA 应用的两个主要难题。PCA 方法在 PLSR 模型和 BP 模型中的表现则大不相同,其是通过降维的思想以及使用较少的综合变量来替代原本较多的变量,且这些综合变量之间相互独立,可以弱化变量自相关所引起的误差,但当光谱中存在较多的非线性信息时,数据降维的效果并不明显<sup>[24-25]</sup>,即在线性的 PLSR 模型中取得的效果较好,而在非线性的 BP 网络中取得的效果一般。SAE 方法通过计算输入光谱数据与输出光谱数据之间的误差,不断调节参数以学习数据内部隐藏的特征,从而压缩输入光谱以提取有用的光谱特征<sup>[26]</sup>。在 SAE 方法的基础上进行 BP 建模所得到的模型精度最高,说明 SAE-BP 方法可以有效地估算 SOM。

从表 2 可以看到,相比于 PLSR 模型下的 SPA、PCA 和 SAE 变量筛选方法,BP 模型下的同种变量筛选方法在测试集上的  $R^2$ 、RPD 和 RPIQ 均有所提高,而 RMSR 有所降低。在 VIS-NIR 光谱中,PLSR 模型仅能处理光谱与 SOM 之间的线性信息而忽略其中的非线性信息,BP 模型作为网络结构具有很强的记忆能力、非线性映射能力以及极强的自学习能力,为此能够学习到高光谱中复杂的非线性特征,从而获得的模型效果更优。

实验采用不同的变量筛选方法,结合 PLSR 模型和 BP 神经网络算法对新疆艾比湖湿地的 SOM 进行反演估算,可以取得较好的效果。然而由于人为的实验误差以及复杂的土壤光谱特性,如何优选光谱以及去除土壤中的干扰因素,建立更稳健的模型仍需进一步研究。

## 5 结 论

在土壤的有机质光谱的建模过程中,光谱数据

的预处理方法以及建模方法都会影响建模的预测效果。以艾比湖湿地的 89 个土壤采样点作为研究对象,采用 SAE、PCA 和 SPA 方法从原始光谱中获取特征变量,建立基于特征变量的 PLSR 和 BP 预测模型。实验结果表明,在土壤的光谱数据中通常包含许多冗余信息,对其进行有效剔除能够提升模型的预测精度。土壤的原始光谱与 SOM 之间的相关性较差,通过对原始光谱进行平滑微分处理能够取得显著效果,其中一阶导数的相关性高于二阶导数;SAE 方法的降维效果优于 SPA 方法和 PCA 方法。对比分析 SAE、PCA 和 SPA 数据降维方法以及 PLSR 模型和 BP 模型,SAE-BP 模型在估算 SOM 中取得的精度最高,得到  $R_2 = 0.82$ 、 $x_{\text{RMSE}} = 3.53$ 、 $x_{\text{RPD}} = 2.32$  和  $x_{\text{RPIQ}} = 2.78$ ,说明该模型能够较好地预测 SOM,这与 SAE 方法和 BP 神经网络都是网络结构以及能够更好地处理光谱中非线性信息相关。

## 参 考 文 献

- [1] Kweon G, Lund E, Maxton C. Soil organic matter and cation-exchange capacity sensing with on-the-go electrical conductivity and optical sensors[J]. *Geoderma*, 2013, 199: 80-89.
- [2] Shi Z, Ji W, Viscarra Rossel R A, et al. Prediction of soil organic matter using a spatially constrained local partial least squares regression and the Chinese VIS-NIR spectral library [J]. *European Journal of Soil Science*, 2015, 66(4): 679-687.
- [3] Wang X P, Zhang F, Kung H T, et al. New methods for improving the remote sensing estimation of soil organic matter content (SOMC) in the Ebinur Lake Wetland National Nature Reserve (ELWNNR) in northwest China[J]. *Remote Sensing of Environment*, 2018, 218: 104-118.
- [4] Ben-Dor E, Banin A. Near-infrared analysis as a rapid method to simultaneously evaluate several soil properties [J]. *Soil Science Society of America Journal*, 1995, 59(2): 364-372.
- [5] Zhang Z P, Ding J L, Wang J Z, et al. Prediction of soil organic matter in northwestern China using fractional-order derivative spectroscopy and modified normalized difference indices [J]. *Catena*, 2020, 185 (10): 104257.
- [6] Li G W, Gao X H, Xiao N W, et al. Estimation of soil organic matter content based on characteristic variable selection and regression methods [J]. *Acta*

- Optica Sinica, 2019, 39(9): 0930002.
- 李冠稳, 高小红, 肖能文, 等. 特征变量选择和回归方法相结合的土壤有机质含量估算[J]. 光学学报, 2019, 39(9): 0930002.
- [7] He D J, Chen X. Real-time measurement of soil organic matter content in field[J]. Transactions of the Chinese Society for Agricultural Machinery, 2015, 46(1): 127-132.  
何东健, 陈煦. 土壤有机质含量田间实时测定方法[J]. 农业机械学报, 2015, 46(1): 127-132.
- [8] Wang H F, Zhang Z T, Karnieli A, et al. Hyperspectral estimation of desert soil organic matter content based on gray correlation-ridge regression model[J]. Transactions of the Chinese Society of Agricultural Engineering, 2018, 34(14): 124-131.  
王海峰, 张智韬, Arnon Karnieli, 等. 基于灰度关联-岭回归的荒漠土壤有机质含量高光谱估算[J]. 农业工程学报, 2018, 34(14): 124-131.
- [9] Shi Z, Wang Q L, Peng J, et al. Development of a national VNIR soil-spectral library for soil classification and prediction of organic matter concentrations[J]. Science China Earth Sciences, 2014, 57(7): 1671-1680.
- [10] Hong Y S, Chen S C, Liu Y L, et al. Combination of fractional order derivative and memory-based learning algorithm to improve the estimation accuracy of soil organic matter by visible and near-infrared spectroscopy[J]. Catena, 2019, 174: 104-116.
- [11] Zhang H L, Luo W, Liu X M, et al. Measurement of soil organic matter with near infrared spectroscopy combined with genetic algorithm and successive projection algorithm[J]. Spectroscopy and Spectral Analysis, 2017, 37(2): 584-587.  
章海亮, 罗微, 刘雪梅, 等. 应用遗传算法结合连续投影算法近红外光谱检测土壤有机质研究[J]. 光谱学与光谱分析, 2017, 37(2): 584-587.
- [12] Luan F M, Zhang X L, Xiong H G, et al. Comparative analysis of soil organic matter content based on different hyperspectral inversion models[J]. Spectroscopy and Spectral Analysis, 2013, 33(1): 196-200.  
栾福明, 张小雷, 熊黑钢, 等. 基于不同模型的土壤有机质含量高光谱反演比较分析[J]. 光谱学与光谱分析, 2013, 33(1): 196-200.
- [13] Ye Q, Jiang X Q, Li X C, et al. Comparison on inversion model of soil organic matter content based on hyperspectral data[J]. Transactions of the Chinese Society for Agricultural Machinery, 2017, 48(3): 164-172.
- 叶勤, 姜雪芹, 李西灿, 等. 基于高光谱数据的土壤有机质含量反演模型比较[J]. 农业机械学报, 2017, 48(3): 164-172.
- [14] Huang F M, Zhang J, Zhou C B, et al. A deep learning algorithm using a fully connected sparse autoencoder neural network for landslide susceptibility prediction[J]. Landslides, 2020, 17(1): 217-229.
- [15] Ayinde B O, Inanc T, Zurada J M. Regularizing deep neural networks by enhancing diversity in feature extraction[J]. IEEE Transactions on Neural Networks and Learning Systems, 2019, 30(9): 2650-2661.
- [16] Feng J, Zhou Z H. Autoencoder by forest [C]// Thirty-Second AAAI Conference on Artificial Intelligence. New York: AAAI, 2018.
- [17] Bello G A, Dawes T J W, Duan J M, et al. Deep-learning cardiac motion analysis for human survival prediction[J]. Nature Machine Intelligence, 2019, 1(2): 95-104.
- [18] Zhang R, Isola P, Efros A A. Split-brain autoencoders: unsupervised learning by cross-channel prediction[C]//2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 21-26, 2017, Honolulu, HI, USA. New York: IEEE, 2017: 645-654.
- [19] Chen B, Zou X Y, Zhu W J. Eliminating outlier samples in near-infrared model by method of PCA-mahalanobis distance[J]. Journal of Jiangsu University (Natural Science Edition), 2008, 29(4): 277-279, 292.  
陈斌, 邹贤勇, 朱文静. PCA结合马氏距离法剔除近红外异常样品[J]. 江苏大学学报(自然科学版), 2008, 29(4): 277-279, 292.
- [20] Ge X Y, Ding J L, Wang J Z, et al. Estimation of soil moisture content based on competitive adaptive reweighted sampling algorithm coupled with machine learning[J]. Acta Optica Sinica, 2018, 38(10): 1030001.  
葛翔宇, 丁建丽, 王敬哲, 等. 基于竞争适应重加权采样算法耦合机器学习的土壤含水量估算[J]. 光学学报, 2018, 38(10): 1030001.
- [21] Hong Y S, Liu Y L, Chen Y Y, et al. Application of fractional-order derivative in the quantitative estimation of soil organic matter content through visible and near-infrared spectroscopy[J]. Geoderma, 2019, 337: 758-769.



- [22] Xu S X, Zhao Y C, Wang M Y, et al. Determination of rice root density from Vis-NIR spectroscopy by support vector machine regression and spectral variable selection techniques[J]. *Catena*, 2017, 157: 12-23.
- [23] Yu L, Hong Y S, Zhou Y, et al. Wavelength variable selection methods for estimation of soil organic matter content using hyperspectral technique[J]. *Transactions of the Chinese Society of Agricultural Engineering*, 2016, 32(13): 95-102.  
于雷, 洪永胜, 周勇, 等. 高光谱估算土壤有机质含量的波长变量筛选方法[J]. *农业工程学报*, 2016, 32(13): 95-102.
- [24] Ouyang Q, Chen Q S, Zhao J W. Intelligent sensing sensory quality of Chinese rice wine using near infrared spectroscopy and nonlinear tools[J]. *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy*, 2016, 154: 42-46.
- [25] Morellos A, Pantazi X E, Moshou D, et al. Machine learning based prediction of soil total nitrogen, organic carbon and moisture content by using VIS-NIR spectroscopy[J]. *Biosystems Engineering*, 2016, 152: 104-116.
- [26] Yuan F N, Zhang L, Shi J T, et al. Theories and applications of auto-encoder neural networks: a literature survey[J]. *Chinese Journal of Computers*, 2019, 42(1): 203-230.  
袁非牛, 章琳, 史劲亭, 等. 自编码神经网络理论及应用综述[J]. *计算机学报*, 2019, 42(1): 203-230.