

Hyperspectral estimation of soil copper concentration based on improved Tabnet model in the eastern junggar coalfield

Yuan Wang, Abdugheni Abliz, Hongbing Ma, Li Liu, Alishir Kurban, Ümüt Halik, Matti Pietikäinen, Wenjuan Wang

Abstract—China is the largest coal consumer in the world. The massive exploitation and utilization of coal resources has resulted in serious problems of heavy metal pollution and environmental contamination, such as soil degradation, water pollution, crop damage, and even threatening human lives. Therefore, monitoring soil heavy metal pollution quickly and in real time is an urgent task at present. This research not only formulated a new preprocessing method enlightened by few-shot learning for soil hyperspectral data, but also combined it with other soil-related auxiliary information to extract effective information from the soil hyperspectrum, at the end of which different regression methods were adopted to predict soil heavy metal contamination. This test used 168 actual soil samples from the Eastern Junggar coalfield in Xinjiang for verification. Since copper in the soil is a trace element and the corresponding spectral characteristics are affected by other impurities, improper use of hyperspectral preprocessing methods may introduce interference information or may delete useful information, which makes the model effect unsatisfied. To effectively address the above problems, the preprocessing method of this experiment combined the second-order differential derivation, data enhancement method together with the

addition of auxiliary information to allow more effective features to be entered into the model. Next, the Attentive Interpretable Tabular Learning (TabNet) model was improved in three different ways using the original TabNet model and three improved TabNet models to create regression models. One of the improved TabNet models had the best effect, with a list of the top 30 features according to the degree of importance. Meanwhile, the regression prediction of Cu content using four different convolutional neural networks (CNN) revealed that the model with the residual block was the strongest and slightly outperformed the improved TabNet model, but lacked interpretation of the input data. Besides, this experiment also employed different pre-processing methods for regression prediction on various models, and found that the traditional pre-processing methods performed best in traditional regression models (e.g., PLSR) and underperformed in deep learning models. The selected optimal model was compared with partial least square regression (PLSR), and convolutional neural network (CNN) models. The results indicated that both the improved TabNet model and improved CNN model had better performance using the new preprocessing approach proposed in this paper, with improved TabNet yielding a coefficient of determination (R^2), root mean square error (RMSE) and ratio of performance to interquartile range (RPIQ) of 0.94, 1.341 and 4.474, respectively. The improved CNN model had a coefficient of determination of 0.942, a root mean square error of 1.324 and an interquartile range of 4.531 in the test dataset.

Manuscript received January 21, 2022; revised March 13, 2022; accepted July 5, 2022. This work was supported in part by the Autonomous Region Postgraduate Research Innovation Project XJ2021G064, in part by National Natural Science Foundation of China 51704259 and in part by Shanghai Aerospace Science and Technology Innovation Fund SAST2019-048. (*Corresponding author: Yuan Wang.*)

Yuan Wang and Abdugheni Abliz contributed equally to this work and should be considered co-first authors.

Yuan Wang is with the College of Information Science and Engineering, Xinjiang University, Urumqi 830046, China, and also with the Department of Electronic Engineering, Tsinghua University, Beijing 100084, China, and also with the Beijing National Research Center for Information Science and Technology (BNRist), Tsinghua University, Beijing 100084 (e-mail: a123penny@163.com).

Abdugheni Abliz is with College of Geography and Remote Sensing Science, Xinjiang University, Urumqi 830046, China (e-mail: abduhini0997@126.com).

Hongbing Ma is with the Department of Electronic Engineering, Tsinghua University, Beijing 100084, China (e-mail: hbma@mail.tsinghua.edu.cn).

Li Liu is with College of System Engineering, National University of Defense Technology, Changsha 410073, China, and also with Center for Machine Vision and Signal Analysis, University of Oulu, Oulu 90014, Finland (e-mail: li.liu@oulu.fi).

Alishir Kurban is with Xinjiang Institute of Ecology and Geography, Chinese Academy of Sciences, Urumqi 830011, China, and also with Sino-Belgian Joint Laboratory for Geo-Information, Xinjiang Institute of Ecology and Geography, Urumqi 830011, China (e-mail: alishir@ms.xjb.ac.cn).

Ümüt Halik is with the school of resources and environmental sciences, Xinjiang University, Urumqi 830046, China (e-mail: umut.halik@qq.com).

Matti Pietikäinen is with the Center for Machine Vision and Signal Analysis, University of Oulu, Oulu 90014, Finland (e-mail: matti.pietikainen@oulu.fi).

Wenjuan Wang is with the Department of Applied Ecology, Saint Petersburg State University, Saint Petersburg 199178, Russian Federation (e-mail: st082527@student.spbu.ru).

Index Terms—Soil hyperspectrum, Soil Cu, Soil heavy metal pollution, Soil auxiliary information, Optimal band combination algorithm, Data enhancement (DA), Few-shot learning, partial least square regression (PLSR), improved convolutional neural network (CNN), improved Attentive Interpretable Tabular Learning (TabNet).

I. INTRODUCTION

AS an important part of the terrestrial ecosystem, soil is an irreplaceable environmental factor for human and animal habitation, as well as a basic guarantee for food safety and human health, and plays an important role in protecting the environment and maintaining ecological balance. With the development of industrial and agricultural production and the urbanization of rural areas, the amount of contaminated soil continues to expand [1]. The large-scale development and utilization of coal resources have made pollution an increasingly serious problem, leading to ecological damage, environmental contamination and other consequences, which in turn affect social and economic development [2], [3]. Therefore, rapid and accurate detection of the copper (Cu) concentration in coal mine areas is of great significance for preventing and ameliorating Cu pollution.

The traditional method of determining heavy metal substances is to collect soil samples in the field and bring them back to the laboratory for chemical testing and analysis, thereby being time-consuming and laborious. With the widespread application of hyperspectral remote sensing in vegetation research, precision agriculture, geological surveys, and military investigations, an increasing number of soil researchers use visible light and near-infrared spectroscopy to determine the heavy metal concentration in soil [4]. However, Cu exists in the soil as a trace element, its reaction in the soil spectrum is relatively weak, and other types of substances, including organic matter, moisture, clay, and iron oxide, in the soil are relatively abundant and thus have an impact on the Cu concentration [5]. The reflectivity in the soil spectrum produces interference, which affects accurate prediction of the Cu concentration. To solve the above problems, Hong-Yan et al. [6] used a method based on first-order reciprocal and baseline correction preprocessing methods to extract the relevant information for Cu analysis. The study found that different preprocessing methods have an effect on Cu determination. Prediction of the elemental concentration has varying degrees of influence. Previous research undertook differential processing of hyperspectral data to reduce the interference of other substances on heavy metal concentration estimation. When measuring the concentration of six heavy metals, Song et al. [7] used three methods, namely, multiplicative scatter correction, unit vector normalization, and the first derivative, to preprocess the soil hyperspectral data before data modeling so that there is more important information in the hyperspectral data, which is convenient for subsequent analysis and detection with the model. Wu et al. [5] explored the relationship between soil hyperspectral data and soil heavy metal concentration. In the preprocessing stage, to reduce the interference of noise, the spectral data were derived using the first and second derivatives, and Savitzky-Golay smoothing was used. In addition, past studies have found that the concentration of other substances in the soil has a significant impact on the relationship between the concentration of heavy metal elements and the soil spectrum [8]. For example, the Cu concentration is related to zinc (Zn), chromium (Cr), arsenic (As), and soil organic matter (SOM). Due to the strong adsorption of organic matter and clay minerals, the heavy metal Ni can be retained in the soil. Sun et al. [9] added the spectral bands related to organic matter and clay minerals as auxiliary information to predict the Ni concentration. To further study whether auxiliary information in soil can help detect the concentration of heavy metals, Hong et al. [10] studied the relationship between the Cr concentration and auxiliary information such as soil organic matter, Fe concentration, pH and the soil spectrum. The study showed that the combination of soil spectral characteristics and auxiliary information was the best way to accurately predict the heavy metal Cr concentration, in contrast to using either of these two broad categories of features individually. Due to the large number of hyperspectral bands, the number of calculations increases during preprocessing, and the effect of model evaluation is not ideal. To reduce the amount of calculation and extract more effective hyperspectral features, previous research adopted spectral feature reduction technology, among which

principle component analysis (PCA) technology and two-band spectral indices are a common method. Yongsheng Hong [10] used the PCA method to extract concise and useful information from the soil spectrum and the optimal band combination algorithm to select the optimal band combination to predict the Cd concentration, and the experimental results show that the optimal band combination algorithm is better than the PCA algorithm.

However, there is often only a small amount of spectral data regarding measured soil, and the traditional regression model used no longer satisfies the precise detection of heavy metal concentration [11], [12]. Besides in the field of soil science, the problem with having just a small number of measured samples can also be found in the field of image [13] and natural language processing [14]. Wang et al. [15] summarized three types of improvement strategies for few-shot learning (FSL), namely, data expansion, model improvement and algorithm improvement. In the image field, many scholars have proposed flipping [16], translation [17], rotation [18] and other strategies to expand data and achieve better results. Spectral information can be regarded as a one-dimensional image signal. Bjerrum et al. expanded spectral data by adding random variations to the offset, multiplication, and slope [19], and repeated the operation nine times. Both the expanded data and the original data together form the training set, in which the expanded data are used to increase the number and diversity of training samples to make the model more robust. Of all the experiments mentioned in this paper, extended multiplicative scatter correction (EMSC) and data augmentation represent the optimal combination of preprocessing and can be used in conjunction with the CNN model to achieve optimal outcomes. In particular, data augmentation seeks to simulate the various types of noise information in the spectrum and, as compared to models that do not utilize expanded data, it can assist convolutional neural networks to concentrate on invariant features, while EMSC allows for an overall baseline correction of the data. To enhance the performance of deep learning models in the field of speech recognition and to prevent the augmentation of data with limited diversity, Wang et al. [20] proposed a generative model which is called a voice conversion technique employing WaveNet (VC-WaveNet) to synthesize speech with diverse pitch patterns to augment data, and demonstrated that the VC-WaveNet technique is an effective data augmentation technique for automatic speech recognition. To surmount the constraints of the small sample size, Wu et al. [21] proposed a new generative adversarial network to generate synthetic Raman spectroscopy data, which was shown to contribute to the accuracy of skin cancer tissue classification, a finding that also demonstrated the reliability of synthetic Raman spectroscopy data. Haut et al. [22] presented that random obscuration in different rectangular spatial regions of hyperspectral images allows the number of samples to increase, which in turn adds to the complexity of the data, and that this effective data augmentation (DA) technique can effectively mitigate the occurrence of overfitting issues, which can result in insufficient generalization and loss of accuracy. The experimental findings from the DA methodology on both hyperspectral datasets indicate that it contributes to better

classification accuracy. However, there is no research on the spectral number expansion method to predict the concentration of heavy metals in soil. The heavy metal concentration estimated by using soil hyperspectral information and traditional regression models is not well studied to organic matter and other nontrace element substances. Whilst models which are less complicated are more general and easier to interpret by researchers, such as PLSR, a significant amount of expertise, time, and experience is required to examine soil hyperspectral data when employing PLSR to predict heavy metal concentration, with the accuracy of the conclusions obtained being somewhat less than the application level. As such, enhancing accuracy by increasing the complexity of the model for predicting heavy metal concentrations could also be an effective breakthrough.

In the field of spectroscopy, deep learning has also been studied more widely. Among them, Zhang et al. [23] predicted corn protein content, tablet active ingredient content, and soil organic carbon content by using an end-to-end deep learning approach, and found that it outperformed the convolutional neural network and three traditional regression methods. Moreover, a deep learning approach combining stacked autoencoders (SAE) and fully connected neural networks (FNN) [24] has been proposed, which has been experimentally demonstrated to be effective in detecting soluble solids concentrations and hardness in the Vis/NIR hyperspectral images of Korla bergamot pear. Furthermore, the research [25] found that hyperspectral image models combined with deep learning-based stacked autoencoders and least squares support vector machines provided the best performance in predicting the total volatile basic nitrogen (TVB-N) concentration of Pacific white shrimp. Liu et al. [26] proposed a deep convolutional neural network for Raman spectra classification that not only offers robust performance, but also eliminates preprocessing steps such as baseline correction or PCA. In the previous research on the prediction of the heavy metal concentration from soil hyperspectral, Pyo et al. [27] performed PCA operations on 98 visible and near-infrared spectroscopy soil samples which resulted in the 1652 bands used being compressed to 68 components. Eventually, the visual and near-infrared spectroscopy was analyzed and processed separately by employing CNN models with convolutional autoencoders, artificial neural network (ANN), and random forest regression (RFR) models. The experiment revealed that CNN models with convolutional autoencoders yielded the highest As, Cu, and lead (Pb) estimates. But the research has a high content and lacks an explanation of the usability of the features. Nevertheless, the low concentration of heavy metals in soil for the our current research increases the difficulty of the present experiment. In order to predict the concentration of compound heavy metals from lettuce leaves, Zhou et al. [28] presented the use of wavelet transform (WT) and stack convolution autoencoder (SCAE) to extract features of compound heavy metals, and ultimately achieved the best prediction outcome with support vector machine regression models. Additionally, some scholars have used data from other mediums in conjunction with deep learning to estimate heavy metal concentration instead of using empirical soil hyperspectral data. Jiang et al. [29] has

used 9 urban multi-source data, which were cost-effective and easily available, as features for input and used an appropriate deep learning algorithm (which is called GRU) to construct predictive models. The model can accurately predict the four heavy metal concentrations of Cu, Zn, Ni, and Cr and is better than the ANN model. Apart from this, a total of two studies exist that use deep learning techniques and soil spectroscopy to predict the concentration of other substances. Qiao et al. [30] employed a technique called SVD concatenation to learn features and combined it with the CNN model (which is called SVD-CNN) to predict the soil organic matters on both FT-NIR and LUCAS 2009 datasets. Ultimately, SVD-CNN was identified as having the highest evaluation and generalization capability. Related studies have also used deep convolutional neural network (DCNN) models to analyze and predict information from soil spectral repositories and have made accurate predictions for most soil properties and outperformed both single-task shallow convolutional neural networks and traditional machine learning methods, demonstrating the potential of modeling with soil spectral data for deep learning [31]. Yet, no proposal has been made to process empirical soil hyperspectral data with data augmentation techniques in few-shot learning and to predict the concentration of heavy metals in soil by applying deep learning models. After the collection and processing of the data used to predict the soil hyperspectral, the spectral data can be stored using CSV tables. Recently, TabNet [32], a deep learning model for tabular data, is proposed, which can discriminate and invert classification and regression tasks. The model proves to have good effects and explanatory properties in many tabular data through experiments.

In this paper, we use the knowledge related to few-shot learning to enhance the soil hyperspectral data, and later modeling with deep learning models, has achieved excellent outcomes. Firstly, to obtain more information from a relatively small number of soil hyperspectral samples, and to tackle the issue that deep learning models are data-driven models, this research introduces a few-shot learning approach that enables deep learning to obtain the Cu concentration in soil from a smaller number of soil hyperspectral samples in a generalized manner. This is accomplished by expanding the hyperspectral data with a priori knowledge, and processing both the expanded data and the original data with second-order derivative for application adding auxiliary information to deep learning models. Next, the TabNet and CNN models, which have shown good performance in handling tabular data, have been used for regression analysis of spectral data to predict the Cu concentration, and by modifying the TabNet and CNN models respectively, the regression models have shown outstanding performance and surpassed the traditional regression models of PLSR. In addition to this, the experiment also lists, illustrates, and analyses the top 30 features chosen by the top-performing TabNet model, effectively revealing the influence of other substances in the soil on the Cu concentration. Last but not least, as far as we know, the present research combines few-shot learning and deep learning for the first time to perform regression analysis on in-situ measured soil hyperspectral data and obtains accurate prediction outcomes. Such an approach

enables the prediction of soil heavy metal concentration to be generalized even when the area is large but the empirical soil hyperspectral data is small, which provides a new concept and methodology for the prediction of heavy metal concentration from soil hyperspectral data and may facilitate its practical application at an early date.

II. MATERIALS AND METHODS

A. Study Area and Sampling Sites

The Eastern Junggar coalfield is located at the north foot of Tianshan Mountain, southeast of the Junggar Basin, within the territory of three counties (Jimsar, Qitai and Muyu) in eastern Changji Prefecture, covering an area of approximately 11,213 km² [33]. This region is located in the hinterland of Eurasia, which has an extremely arid and continental climate. Winter is long and cold, and summer is short and hot. The annual average temperature is 7 °C, and the annual average precipitation is 183.5 mm. The difference between spring and autumn is not obvious. The soil types in Zhundong coal mining area mainly consist of saline soil, eolian sandy soil, gray-brown desert soil, gypsum brown desert soil, and desert alkali soil [33]. In this study the soil samples taken from the same soil type called gray-brown desert soil which was located around the coal mining region. The plant composition is simple and monotonous, mainly xerophytic and super-xerophytic shrubs, semi-trees and herbs, including Tamarix, Haloxylon ammodendron, Pilosa and other plants.

Depending on the topographic feature of the study area, as centered of coal mining region and chemical plant, considering the direction of pollutant emission from industrial areas (mainly northwest wind in this region), therefore, soil samples were collected from field by using systematic random sampling method. A total of 168 soil samples were systematic randomly collected from June to July 2013, and a 5-point mixed sampling method (0-20 cm deep) was adopted for each soil sample. Due to the characteristics of the Eastern Junggar open-pit coal mine, a handheld GPS device was used to locate the soil sample points, as shown in Fig. 1. The collected soil was packaged in a plastic bag, brought back to the laboratory, dried, ground, and passed through a 0.2 mm aperture sieve. Part of the processed soil samples was collected for soil hyperspectral data, and the other part was entrusted to the Physical and Chemical Testing Center of Xinjiang University for determination by potassium dichromate-volumetric dilution calorimetry [34].

B. Laboratory Spectral Measurements and Preprocessing

The soil spectrum was measured by an ASD Field Spec3 portable spectroradiometer (Analytical Spectral Devices, Inc., USA), and its spectral range was 350-2500 nm. The spectrum measurement was collected in a dark room to avoid the influence of external light. The light source employed was a halogen lamp with a power of 50 W, which was used for the reflection spectrum of the sample. The vertical angle was 25°, 0.5 meters away from the sample. The field of view of the spectrometer probe was 25°, perpendicular to the soil sample, and 15 cm away from the sample. A spectralon white panel

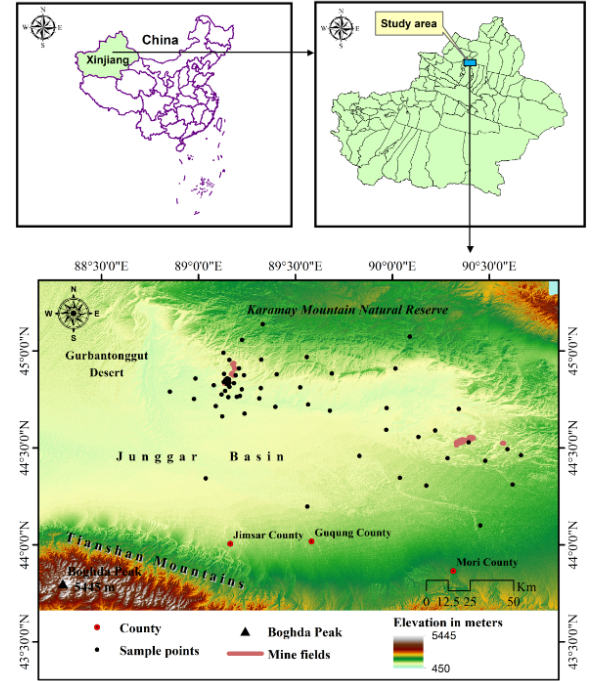


Fig. 1. Study areas and location of sampling points.

was applied as the standard reference for diffuse reflection, and it needed to be calibrated at the same time as the measurement. Each sample was measured 10 times. If there was an abnormal spectrum, it was eliminated, and the arithmetic average of the 10 measurements was taken as the spectral reflectance of the sample.

To reduce random noise, the data at both ends of the spectrum (350-399 and 2401-2500 nm) were removed. A Savitzky-Golay filter was employed to smooth the 401-2400 nm spectral curve [35]. Origin 8.0 software was used for spectral smoothing and noise removal processing. To eliminate or weaken the background noise, enhance the target spectral information, and improve the signal-to-noise ratio, the second-order reciprocal [36] was applied to preprocess the spectral data. Finally, the spectrum was resampled with a resampling interval of 10 nm. In this preprocessing task, we used a personal computer with a 2.8 GHz Intel Core i7-4900 MQ, 32 GB RAM, and a Windows 10 operating system. Compared to that of traditional methods, our preprocessing time was less than 20 minutes, which saves time.

C. FOD

Fractional order differentiation can provide an extension of the order of integer order differentiation to arbitrary orders [37]. Currently, the three main forms that are commonly used are Riemann-Liouville, Grünwald-Letnikov and Caputo [38]. Of these, the Grünwald-Letnikov expression is defined as:

$$d^\alpha f(x) = \lim_{h \rightarrow 0} \frac{1}{h^\alpha} \sum_{m=0}^{\frac{t-a}{h}} (-1)^m \frac{\Gamma(\alpha+1)}{m! \Gamma(\alpha-m+1)} f(x-mh) \quad (1)$$

In which α is an arbitrary order, h is the step size, while t and a are the maximum and minimum limits of the differentiation, respectively.

Given the resampling interval of 1 nm for the spectrometer used in the research, with $h = 1$, it is possible to derive the difference expression for the fractional order differentiation of the unary function $f(x)$ as:

$$\begin{aligned} \frac{d^\alpha f(x)}{dx^\alpha} \approx & f(x) + (-\alpha)f(x-1) + \frac{(-\alpha)(-\alpha+1)}{2}f(x-2) \\ & + \cdots + \frac{\Gamma(-\alpha+1)}{n!\Gamma(-\alpha+n+1)}f(x-n) \end{aligned} \quad (2)$$

D. Data Enhancement

Deep learning is a data-driven method using multilayer neural networks to learn different levels of detailed features, which is more conducive to the improvement of regression model performance. In the existing spectroscopy research, due to the small number of datasets, the spectroscopy expands through stochastic feature mapping [39], cross domain feature adaptation [40], Generative Adversarial Network (GAN) [41], etc. to make the model robust and adaptable. Compared with two-dimensional data, such as images, the spectrum can be regarded as one-dimensional data. Esben et al. [19] uses offset, slope, and multiplication spectral data expansion methods to expand the tablet spectral data, and combined with deep convolutional neural networks, it has a strong regression effect. In the experiments mentioned above, multiplication was done 1 ± 0.10 times, but the tablet dataset from near infrared (NIR) spectra was used to measure drug content, and the spectral response was related to drug content. For this experiment, on the other hand, the soil hyperspectral response is related to the Cu concentration, which is why the choice of expanded parameters and the above experiments should be considered separately in conjunction with specialized domain knowledge. As the present experiment uses second-order differential to preprocess hyperspectral data, it is effective and reasonable to combine the enhancement process with the preprocessing method in consideration. In previous spectral preprocessing studies, the fractional-order differential preprocessing method has been used to extract more subtle features with positive outcomes. For the parameters chosen in the present research, the disturbed spectral values range as far as possible between the spectra processed by second-order differential and 2.2 order differential, and between the spectra processed by second-order differential and 1.8 order differential, so as to make the disturbance more reasonable and effective. Following the second-order differential and 1.8 order differential for all bands, the absolute value of the difference for each band lies mostly in the range of 0-0.001 multiplied by the band after the second-order differential, while following the same operation for all bands after the second-order differential and 2.2 order differential, the absolute value lies mostly in the range of 0-0.001 multiplied by the band after second-order differential (Fig. 2). Hence, this research chose to randomly multiply by 0.999 -1.001 for all the soil spectrum datasets.

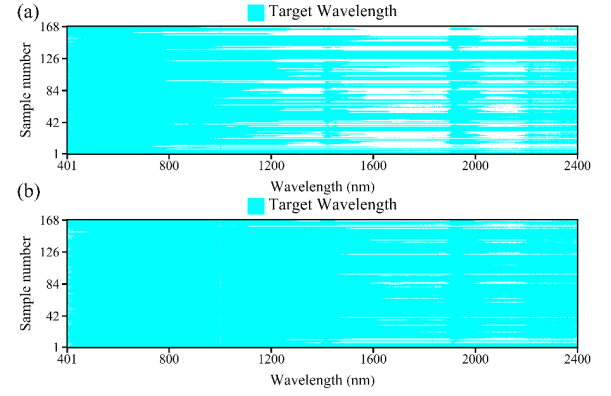


Fig. 2. Distribution of target wavelengths based on the absolute value of the difference in the range of 0-0.001 multiplied by the band after second-order differential (400-2400 nm). The color blue represents the corresponding selected wavelength. (a) the absolute value of the difference between the spectra processed by second-order differential and 2.2 order differential, (b) the absolute value of the difference between the spectra processed by second-order differential and 1.8 order differential.

There are many advantages of using data enhancement as the preprocessing method. Firstly, it is a type of few-shot learning, many expansion methods are already available in other spectral domains, while there is no relevant expansion method for predicting heavy metal concentration from soil hyperspectral measurements, and it has much potential to be combined with more expert knowledge to generate more effective expansion methods. Secondly, the expanded data can effectively alleviate the underfitting issue caused by the insufficient sample size of the deep learning dataset, enabling the deep learning to perform more effectively. Moreover, by performing operations on a larger range of data enables the model to concentrate on statistically significant factors, which can reduce potential bias. Ultimately, the combination of deep learning data augmentation and expanded data from soil hyperspectral heavy metal knowledge eliminates the need for extensive manpower and resources to collect and analyse massive datasets, which reduces expenses and provides cost savings.

E. Spectral Feature Reduction

In recent years, the dual-band spectral index method has been repeatedly employed to explore the relationship between soil heavy metals and wavebands. In this experiment, two dual-band spectral indices were tested for correlation with soil Cu concentration, and calculated using excel sheet software for any two wavelengths (S_i, S_j) within 400-2400 nm. For either of these two spectral indices, the combination of wavelengths with the greatest correlation to Cu concentration was selected using the R software [42]. Normalized difference index (NDI), ratio index (RI) are defined by:

$$NDI(S_i, S_j) = (S_i - S_j) / (S_i + S_j) \quad (3)$$

$$RI(S_i, S_j) = S_i / S_j \quad (4)$$

where S_i and S_j represent the spectral values, and i and j stand for any band from 400 nm to 2400 nm.

PCA is an effective technique to reduce the dimensionality of the data. Through the orthogonal process, hyperspectral data can be transformed into uncorrelated principal components (PCs), which can be linearly combined to represent the original spectral data [43]. Typically, a combination of some well-ranked PCs together can represent most of the raw spectral information. In this study, the number of PCs were based on the explained amount of total spectral variance ($>95\%$) of the total spectral dataset. We used the PCA method in our experiments on both the raw hyperspectral data and the expanded hyperspectral data, all of which were processed by second-order differentiation.

F. Predictive Mechanism Exploration

Many existing studies have shown that other substances in the soil can affect the response of the spectral curve and thus have an impact on the predicted heavy metal concentration in the soil [44]. Given the complex composition of soil and the different concentrations of substances in different regions, the bands and intensity of the influence of these substances on the hyperspectral data vary, thereby affecting the experimental results using spectroscopy to detect soil heavy metal levels to varying degrees. For example, Somsubhra used VisNIR DRS to quickly predict the total arsenic concentration in five different solid arsenic items. After predicting the hyperspectral spectrum and experimenting with the PLSR model, a good result was finally obtained. Through qualitative spectral analysis and PLSR coefficients, the predicted As concentration was shown to have a high correlation with soil organic matter, clay minerals, Fe and Al oxide [45]. Khosravi et al. [46] used hyperspectral data in the 350-2500 nm range for preprocessing using Savitzky-Golay (SG) smoothing, first derivative (FD), and second derivative (SD) approaches and then used partial least square regression (PLSR) and extreme learning machine (ELM) to predict the heavy metal concentrations of Pb and Zn in the soil. Among the algorithms, FD-ELM had the highest prediction accuracy. This result shows that the adsorption of heavy metals by active iron oxides and clay in the soil is an important mechanism for predicting Pb and Zn levels without spectral characteristics. Therefore, this experiment uses Cu and SOM, As, Cr, Pb, and Zn to perform correlation analysis, and it is found that the wavelengths related to the hyperspectral data of the Eastern Junggar coalfield have an analyzable distribution. When selecting training data, auxiliary substances such as SOM, As, Cr, Pb, and Zn can be added so that these substances can be used to replace spectral features when the spectral response is weak.

III. CORRECTION STRATEGY

A. PLSR

Partial least square regression analysis is an effective multiple regression modeling method that has been widely used to establish spectral quantitative models [47]. If the independent variables have serious multiple correlations, the use of stepwise regression to select variables increases the interpretable error of the model. To solve the interference of multiple correlations of variables on regression modeling, Wood, Abano

and others proposed a partial least square regression analysis method, which dealt with a small number of samples, serious multiple correlations between variables, and a large number of explanatory variables. This method has unique advantages, and it can realize regression modeling, data structure simplification and correlation analysis between two sets of variables. Regarding the hyperspectral determination of heavy metals, the spectral absorption characteristics of heavy metals in soil and the influence of certain soil components can be used to predict the concentration of heavy metals in the soil using reflectance spectra. Due to the high correlation between the multiple hyperspectral bands, the partial least square method can solve the above problems to a large extent, and the partial least square regression analysis method is widely used in prediction of the soil heavy metal concentration.

B. CNN

A convolutional neural network [48] is a kind of feed-forward neural network with a convolutional layer and deep network structure, and it is one of the representative algorithms of deep learning. In the convolutional neural network, the convolutional layer and the pool sampling layer of the hidden layer are the core modules to realize the feature extraction function of the convolutional neural network. The network model uses the gradient descent method to minimize the loss function to reversely adjust the weight parameters in the network layer by layer and improves the accuracy of the network through iterative training. Convolutional neural networks not only have a high-quality performance in two-dimensional data but also demonstrate a strong performance in one-dimensional data. For one-dimensional hyperspectral data, each convolution filter in the convolutional layer convolves the input spectral information, and the convolution result constitutes a feature map of the input spectral signal. Different convolution kernels extract different levels of features of the spectral signal and uses these features to improve the final prediction accuracy. Each convolution filter shares the same parameters, including the same weight matrix and bias terms. The benefit of sharing weights is that the location of local features does not need to be considered when extracting features from high-dimensional spectral signals, and weight sharing provides an effective way to greatly reduce the number of convolutional neural network model parameters to be learned. To enhance the model's ability to fit nonlinear measured data, a convolutional neural network is added to the relu' layer. The ReLU layer generates the output of some neurons 0, making the neural network sparse and reducing the interdependence between parameters, thereby decreasing the occurrence of overfitting. The fully connected layer is also often used in convolutional neural networks. The fully connected layer uses all of the local features to improve the models learnability.

C. TabNet

The TabNet model [32] is a high-accuracy and interpretable deep learning model specially designed for tabular data. Not only does it use an end-to-end model to reduce data processing time, but it also uses sequential attention to select features,

so that the reasons for feature selection at each step can be understood, enabling the model to be better interpretable. The feature selection in the model is performed using the attentive transformer layer, which gives the Mask matrix of the current step based on the results of the previous step and attempts to ensure that the Mask matrix is as sparse and non-repetitive as possible. Because the mask vector of the sample can be different, TabNet can let different samples choose various features. The feature transformer layer realizes the calculation and processing of the features selected in the current step. The decision tree structure is a combination of the size relationship of a single feature, which is the decision manifold. TabNet performs feature calculation through a more complex feature transformer layer and is more efficient than decision trees in feature combination.

TabNet enables the use of a tree-like function to obtain a specific value to measure the importance of a feature. TabNet is based on a tree-like function that determines the proportion of each feature through the composition factor Mask. To begin with, TabNet constructs a sequential multi-step structure, where at each step of the decision tree, the most salient features from each decision step are chosen for downstream tasks by using a sequential attention mechanism, making it interpretable and enabling better learning. In addition, TabNet employs a single deep learning framework for end-to-end learning. Secondly, TabNet outperforms, or is on par with, other tabular learning methodologies. It has two types of interpretability, a local interpretative, which demonstrates the importance of each input feature and how they are combined, and a global interpretative, which quantifies the contribution of each input feature in the output.

In this case, the same dimensional features $f \in R^{B \times D}$ are used to each decision step, feature selection is achieved by means of the attentive transformer. The learnable mask $M[i] \in R^{B \times D}$, which is used as a soft selection of salient features from i^{th} step, where B is the batch size and D is the dimensional features. With sparse selection of most salient features, the ability to learn irrelevant features at each decision step would not be wasted. The mask adopts the form of product, that is, $M[i] * f$, based on the feature information processed in the previous step $a[i-1]$, and uses the attentive transformer back to the mask as follows:

$$M[i] = \text{sparsemax}(P[i-1] \cdot h_i(a[i-1])) \quad (5)$$

Of which,

$$\sum_{j=1}^D M[i]_{b,j} = 1 \quad (6)$$

where h_i is a trainable function, as in the part of the attentive transformer above, including a fully-connected (FC) layer, followed by batch normalization (BN). $P[i]$ refers to the prior scale of the above graph. This indicates how many features have been used previously.

$$P[i] = \prod_{j=1}^i (\gamma - M[j]) \quad (7)$$

where the γ in the above equation serves as a relaxation parameter. A feature is used in only one decision step when

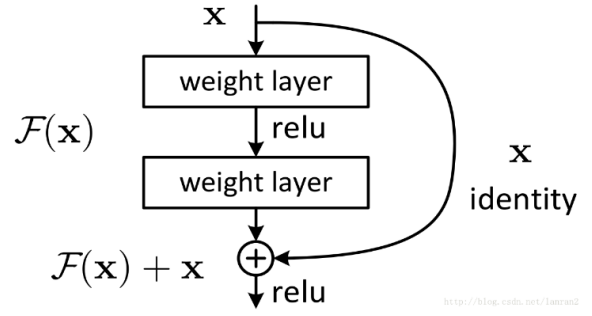


Fig. 3. After formally denoting the desired underlying mapping as $H(x)$, we let the stacked nonlinear layers fit another mapping of $F(x) := H(x) - x$. The original mapping is recast into $F(x) + x$.

$\gamma = 1$; it can be used in multiple decision steps when γ increases.

The mask of each step represents a local interpretation, while the converged mask represents a global interpretation, where the contribution of each input to the output can be obtained by aggregating feature importance mask of all steps:

$$M_{agg-b,j} = \sum_{i=1}^{N_{steps}} \eta_b[i] M_{b,j}[i] / \sum_{j=1}^D \sum_{i=1}^{N_{steps}} \eta_b[i] M_{b,j}[i] \quad (8)$$

where $\eta_b[i]$ denotes the aggregate decision contribution at i^{th} decision step for the b^{th} sample.

D. Residual Learning

As the network deepens, a drop in the accuracy of the training set occurs and is not caused by overfitting. Therefore, the author proposes a brand-new network for this problem, called a deep residual network [49] to allow the network to deepen as much as possible. The residual block is an essential part of the deep residual learning framework. The structure of the residual block is shown in Fig. 3. It does a reference (x) for the input of each layer and learns to form a residual function instead of learning some functions without reference (x). This residual function is easier to optimize and can greatly deepen the number of network layers.

E. Prediction Accuracy

In the process of training set verification and test set verification, this experiment uses commonly used indicators, namely, the coefficient of determination (R^2), root mean square error (RMSE) and performance ratio interquartile range (RPIQ), to evaluate the performance of the model [50]. In the RPIQ evaluation [51], $2.20 \leq \text{RPIQ} < 2.70$ indicates that the model performance is very poor, $2.70 \leq \text{RPIQ} < 3.37$ indicates that the model performance is average and can provide reasonable estimation results, $3.37 \leq \text{RPIQ} < 4.05$ indicates that the model has good performance, and $\text{RPIQ} \geq 4.05$ shows that the model performance is considered excellent. In most cases, good models are considered to have relatively large R^2 and RPIQ values and relatively low RMSE values.

TABLE I
STATISTICAL DESCRIPTIONS OF SOIL CU AND OTHER SOIL PROPERTIES

Attributes	Sample sets	N	Min mg kg ⁻¹	Max mg kg ⁻¹	Mean mg kg ⁻¹	Median mg kg ⁻¹	Standard deviation	CV	Background value
Cu	Entire	168	0.0110	0.0460	0.0192	0.0185	0.0055	0.2866	26.7
	Training	1340	0.0120	0.0460	0.0190	0.0181	0.0050	0.2659	26.7
	Validation	170	0.0110	0.0320	0.0202	0.0223	0.0067	0.3313	26.7
	Test	17	0.0080	0.0340	0.0199	0.0193	0.0075	0.3800	26.7
SOM	Entire	168	0.2634	95.9048	6.3468	2.9289	10.1851	1.6047	None
As	Entire	168	1.5427	80.4553	37.9147	37.3292	6.8191	0.1798	11.2
Cr	Entire	168	0.0140	0.1100	0.0536	0.0540	0.0183	0.3423	49.3
Pb	Entire	168	0.0040	0.0470	0.0214	0.0210	0.0076	0.3544	19.4
Zn	Entire	168	0.0168	0.1090	0.0476	0.0476	0.0144	0.3038	68.8

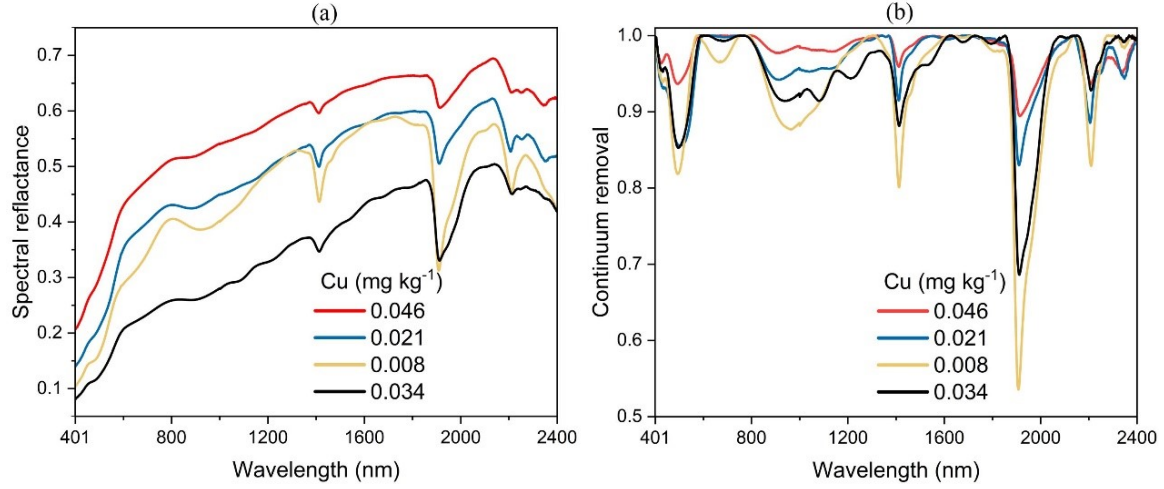


Fig. 4. (a) Original spectral reflectance of soil samples with four different Cu contamination levels and (b) continuum removal (CR) spectra of soil samples with four different Cu contamination levels.

IV. RESULTS

A. Descriptive Statistics for Cd and Other Soil Properties

The statistical data set for the entire Cu concentration is shown in Table I, the minimum value is 0.0110, the maximum value is 0.0460, the mean value is 0.0192, and the coefficient of variation (CV) is 28.66%. According to the pollution level in the China Soil Environmental Quality Control Standard [10], only 39 of the 168 samples had Cu pollution problems, and 129 samples were not contaminated, with a pollution rate of 23.21%.

The SOM, As, Cr, Pb and Zn values ranged from 0.2634 to 95.9048, 1.5427 to 80.4553, 0.0140 to 0.1100, 0.0040 to 0.0470, 0.0168 to 0.1090, respectively (Table I). According to the results reported by Wilding [51], datasets with $CV > 35\%$, $15\% < CV < 35\%$ and $CV < 15\%$ were considered to have high, moderate and low variability, respectively. The CV values of As, Cr and Zn belonged to moderate variation, while the CV value of SOM and Pb were high.

B. Soil Spectral Analysis

Fig. 4 shows that the spectral reflectance curves of soil samples with different Cu concentrations in the study area are roughly the same, and the spectral reflectance of soil samples

and their copper concentration are no longer simply linearly correlated. The reflectance rises rapidly in the visible light range and tends to increase slowly in the near-infrared band. The spectrum curve has 3 more obvious moisture absorption peaks near 1413 nm, 1922 nm, and 2200 nm caused by water molecules and -OH groups. The iron oxides absorbance highly relative to the reflectance in infrared-near infrared band [52], therefore, it would be main influencing factor of causing the reflectance differences in 800-1000 nm. SOM is the most important adsorbent for metals; thus, it was considered to be important factor to determining the species and bioavailability soil heavy metals [53], [54]. Therefore, SOM would be another main influencing factor of causing the reflectance differences in 800-1000 nm. Most of the Cu concentration soil spectra in the range of 401 nm to 2400 nm overlap, and most of the spectra have the same trend, making it difficult to distinguish. Therefore, the spectra are pre-processed using second-order differentiation, and deep learning is applied to synthesize features at each granularity to make the optimal prediction.

C. Correlations Between Cd and Optimal Spectral Indices

The current experiment achieved optimal outcomes with a modified TabNet model. Moreover, the experiment also

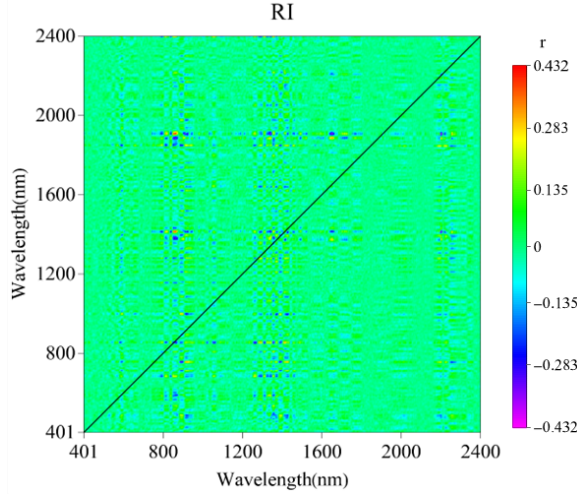


Fig. 5. Correlations (r) between soil Cu and ratio index (RI).

discovered several bands that can have a significant impact on the prediction of Cu concentration and revealed the top 30 bands used in the model. Of these, the bands associated with organic matter are the most numerous, containing a total of 12 bands. Among these, bands primarily correlated with organic C-H are 1761 nm [55], 2380 nm [56], 1829 nm [57], 2320 nm [58], 2327 nm [55], 2321 nm [58], and 2365 nm [58]. The two bands of 1947 nm [56], and 2036 nm [58], [59] are more clearly influenced by C=O in the organic matter, while the chlorophyll pigment and CH₂, CH₃ in the organic matter lead to fluctuations at 631 nm [55], 2284 nm [55], [59], [60] respectively. In the Fe oxides category, the occurrence of Fe-OH, gibbsite, ferrihydrite, water molecular vibrations, and OH⁻ are associated with bands of 897 nm [61], [62], 2266 nm [59], 972 nm [63], [64], and 1422 nm [59], [65] respectively. Among the clay minerals, 1985 nm [57], [66], 1909 nm [57], [67], and 2152 nm [58], [59] are identified as the significant bands for Cu concentration detection. The carbonates, the hydroxyl group, and the Al-OH lattice structure, Al-OH, and kaolin contained in the soil perform an essential role in these three bands respectively.

For this research, the Pearson correlation (r) between Cu concentration in soil and the RI and NDI indices were calculated separately by employing the optimal band combination algorithm. For SDR, the highest correlation regions of NDI and RI were mainly located over the range of 1800-2400 nm (Figs. 5 and 6). Considering the NDI indices for all spectra, which were selected by the degree of correlation and by using the top four spectral indices ranked by the NDI strategy, resulted in NDI ($r = -0.4378$, NDI [SDR 623, SDR 2091]), NDI ($r = 0.4343$, NDI [SDR 851, SDR 880]), NDI ($r = 0.4581$, NDI [SDR 1581, SDR 2331]) and NDI ($r = 0.4571$, NDI [SDR 1734, SDR 2143]). Overall, after the FDR transformation, the four highest Pearson correlation spectral indices of RI were RI (SDR 403, SDR 2195), RI (SDR 870, SDR 2385), RI (SDR 688, SDR 1851), and RI (SDR 1949, SDR 2195), with the Pearson correlation value of 0.4298, 0.4242, -0.4267, and -0.4301, respectively.

In this experiment, owing to the richness of information in

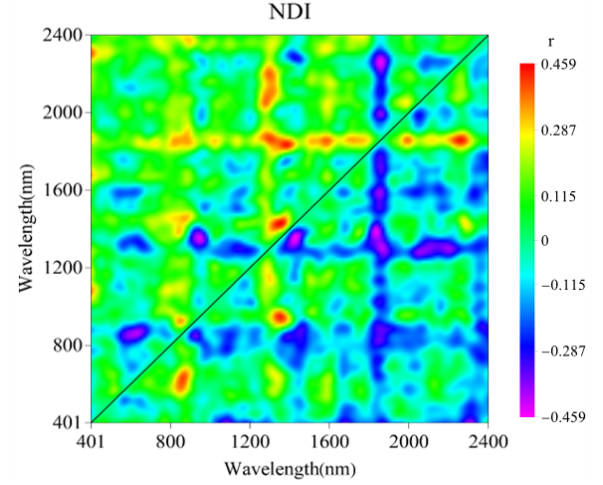


Fig. 6. Correlations (r) between soil Cu and normalized difference index (NDI).

the expanded and unexpanded spectral data used, the required number of Principal components was also higher when the cumulative variances were at nearly 95% with smaller values of the largest explained variances (Fig. 7). The diversity of information adds difficulty to the prediction of Cu content. Although the use of PCA can extract valid information and reduce unnecessary noise information, it still limits the high accuracy of the PLSR implementation. While comparing pre-processing methods that incorporate expanded data yields the best results (Table IV, Table VI), it still has inferior accuracy compared to the combination of also using PCA methods and PLSR models [10]. Furthermore, the use of PCA algorithms results in dimensionality reduction losses, and the relatively small amount of information can cause underfitting of deep learning models, resulting in poorly performing deep learning network models (Table VIII).

D. Comparisons Between the TabNet Model and the Three Improved TabNet Models

Our experiment is based on the TabNet model and is improved. To test the effect of the improved TabNet model, we changed the TabNet model to the FC layer in three different ways (Fig. 8), used the original TabNet [see Fig. 8(a)] to experiment with soil hyperspectral data, and finally compared the four models final results. First, when inputting data, the data of two-band (2D) indices, and auxiliary information, the soil spectrum data processed by the second-order differentiation are used. Spectral concentration is input, and self-supervised learning is used to learn the relationship between features. This experimental procedure is the same for the four experiments. Second, all models use the end-to-end model, where the encoder architecture is shown in Fig. 9, and the TabNet decoder architecture is shown in Fig. 10. The four model transformation parts are the feature transformer parts. Model 1 completely migrates the feature transformer part [see Fig. 8(a)], and to make the input part of the model more abundant, Model 2 adds an FC layer to the shared across decision steps section [see Fig. 8(b)] compared with Model

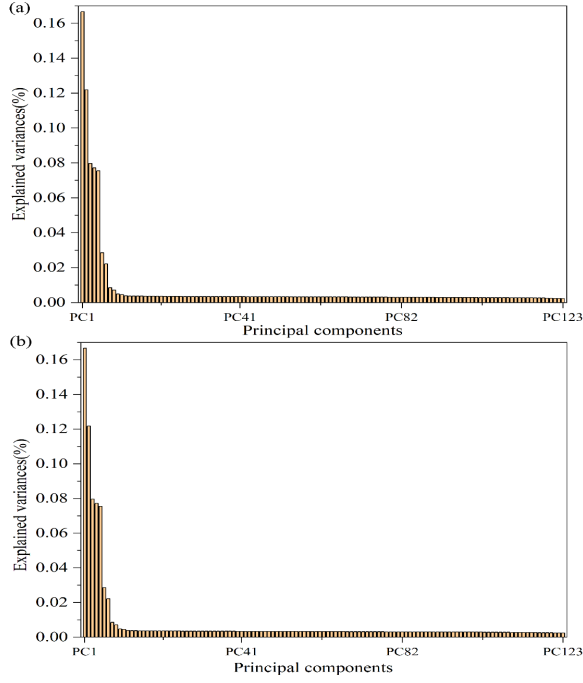


Fig. 7. Correlations of soil Cu to the first 123 principal components. (a) calibration dataset with $n = 1340$. (b) calibration dataset with $n = 134$.

1. Model 3 adds 2 FC layers to the shared across decision steps part [see Fig. 8(c)] on the basis of Model 1, which is used to test whether the more fully connected layers are input can create a better effect. Model 4 not only adds the FC layer on the basis of Model 1 but also adds a residual network to the FC layer for the fully connected layer [see Fig. 8(d)] to reduce information loss when transmitting information and protecting the integrity of the information. It can be seen from Table II that when the number of FC layers increases by 1 layer, the RMSE value of the training set and the validation set increases, and the R^2 value and the RPIQ value decrease. While the RMSE value of the test set decreases, the R^2 value and the RPIQ value increases, thus improving the models effect. When the number of FC layers increases by 2 layers, the RMSE values of the training set and the validation set are increased from 3.07, 3.32 to 3.26, 3.42, respectively, the R^2 value is reduced from 0.6818, 0.6839 to 0.6397, and 0.6656, respectively, the RPIQ value is reduced from 1.95, 2.10 to 1.83, and 2.04, respectively. In the test set, the RMSE value increases from 3.45 to 3.73, the R^2 value decreases from 0.60 to 0.53, and the RPIQ value decreases from 1.73 to 1.60. At this time, when the FC value increases, the model performance is overfitted, and the effect deteriorates. In model four, after adding the residual network, the training set, validation set, and test set have the smallest RMSE value, the largest R^2 value and RPIQ value, and the performances of these four models are the best, far exceeding other models.

E. Impact of Change in Sample Size on the Models Performance

The experiment explores the influence of the changes in the number of training sessions on the results of the experiment.

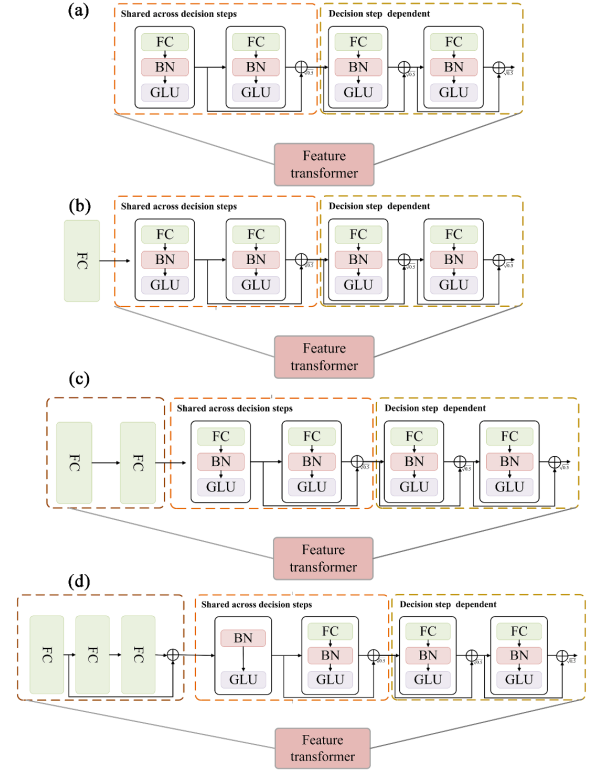


Fig. 8. A feature transformer block example. (a) Model 1, (b) Model 2, (c) Model 3, (d) Model 4.

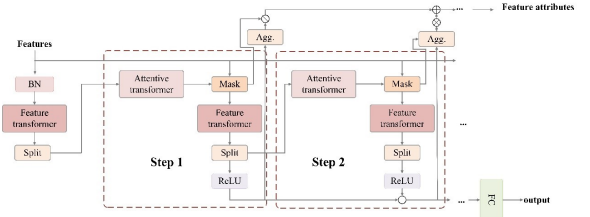


Fig. 9. TabNet encoder, was comprised of a feature transformer, an attentive transformer and feature masking.

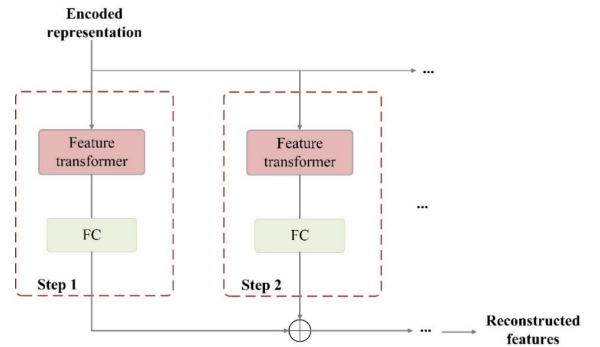


Fig. 10. TabNet decoder, was comprised of a feature transformer block at each step.

For Model 4, which has the best model effect, half of the dataset is used for training, and the results are obtained. After using the TabNet method, it can be seen from Table III that when the amount of training data increases, the RMSE value

TABLE II

STATISTICAL SUMMARY OF CU ESTIMATION WITH THE PRETREATMENT OF STRATEGY I DEVELOPED FROM MODEL 1, 2, 3 AND 4. STRATEGY I INDICATES THAT THE TYPES OF INPUT VARIABLES CONTAINS BOTH THE EXPANDED DATA AND ORIGINAL SPECTROSCOPY DATA WITH SECOND-ORDER DERIVATIVE ADDING AUXILIARY INFORMATION

Strategies	Models	Calibration ($n = 1340$)			Valid ($n = 17$)			test ($n = 17$)		
		r^2	RMSE mg kg ⁻¹	RPIQ	r^2	RMSE mg kg ⁻¹	RPIQ	r^2	RMSE mg kg ⁻¹	RPIQ
Strategy I	Model 1	0.762	2.656	2.259	0.817	2.528	2.769	0.578	3.570	1.680
Strategy I	Model 2	0.682	3.072	1.953	0.684	3.326	2.104	0.606	3.450	1.739
Strategy I	Model 3	0.640	3.269	1.835	0.666	3.422	2.046	0.537	3.738	1.605
Strategy I	Model 4	0.967	0.992	6.051	0.970	1.022	6.852	0.940	1.341	4.474

TABLE III

EFFECT OF THE NUMBER OF CALIBRATION SAMPLE SIZE WITH THE PRETREATMENT OF STRATEGY I ON THE EVALUATION CRITERIA VALUES. STRATEGY I INDICATES THAT THE TYPES OF INPUT VARIABLES CONTAINS BOTH THE EXPANDED DATA AND ORIGINAL SPECTROSCOPY DATA WITH SECOND-ORDER DERIVATIVE ADDING AUXILIARY INFORMATION

Strategy	Strategy I								
	Calibration ($n = 1340$)			Valid ($n = 17$)			test ($n = 17$)		
	r^2	RMSE mg kg ⁻¹	RPIQ	r^2	RMSE mg kg ⁻¹	RPIQ	r^2	RMSE mg kg ⁻¹	RPIQ
Calibration sample size									
13400	0.967	0.992	6.051	0.970	1.022	6.852	0.940	1.341	4.475
6700	0.944	1.278	4.695	0.941	1.433	4.886	0.912	1.630	3.681

decreases from 1.63 to 1.34, which shows that increasing the dataset can increase the model's ability to learn from the data. Since the TabNet model uses different runs on the same dataset, there are different results. Table III provides the average and standard deviation of the RMSEP for 20 runs of the dataset. Although the hyperparameters are fixed, the weights after training will vary as the model runs. This is due to the influence of some random factors, such as random loss of weights and random optimization. In neural networks, the most common way to use randomness is the random initialization of network weights, such as the randomness of weight regularization, the randomness of the dropout layer, and random optimization. In addition to these, there are also even more sources of randomness, meaning that when running the same neural network algorithm on the same dataset, it is destined to receive different results.

F. Three PLSR Models and Three Pretreatments

First, the traditional algorithm PLSR analyses all bands with Strategy I and obtains regression results. In partial least square modeling, the number of components is extracted. In most cases, the partial least square method decomposes both the dependent variable and the independent variable, arranges the factors according to the correlation between them from large to small, and decides to choose a certain number of factors to participate in the modeling. The choice of the number of factors is determined experimentally, which has a highly important impact on the results of the experiment. In this experiment, the performance of the model is evaluated by changing the number of principal components. The parameters and results in the experiment are shown in Table IV. When the number of LV is 5, the model performance is the best, and when the number of LV is 2, the model effect is the worst. This is because when the principal component number is small, the data cannot be

effectively decomposed and filtered, and the comprehensive variable that has the strongest explanation for the dependent variable can be extracted. When the number of principal components gradually increases, the RMSE of the training set and the validation set gradually decreases, R^2 gradually increases, and RPIQ gradually increases, indicating that when more principal components are selected, the model can obtain more Cu concentration information from the hyperspectral soil data. In the test set, when the selected number of principal components gradually increases, the RMSE value, absolute value of R^2 , and RPIQ first decrease and then increase, indicating that when the number of principal components increases, more information is extracted from the hyperspectral spectrum, but it is not useful information related to the Cu concentration. As the information gradually increases, some irrelevant information is introduced, which does not improve the performance of the model and even has a negative impact. Therefore, the performance of the model improves and then worsens. When the number of factors involved in modeling is 5, the result is optimal, and the RMSE, R^2 , and RPIQ of the test set are 5.8104, -0.1179, and 1.0326, respectively.

When applying the pretreatment of Strategy II the PLSR model outperformed the augmented data effect (Table IV, Table V). When using the spectral and auxiliary data after second-order differentiation (Table V), the best result obtained was when the number of LVs was 2, with R^2 values of -0.186 for the test dataset, 5.458 for RMSE and 1.099 for RPIQ. This optimal result outperformed the strategy I preprocessing method, yet the R^2 for both preprocessing methods was negative, illustrating that when the number of training sets was declined, the interference by noise on the PLSR model could be mitigated, albeit the performance of the model remained relatively poor. As the LV value increased, the model became progressively less effective. Whereas PLSR showed the best

TABLE IV

EFFECT OF THE LV NUMBER AND THE PRETREATMENT OF STRATEGY I ON THE EVALUATION CRITERIA VALUES. STRATEGY I INDICATES THAT THE TYPES OF INPUT VARIABLES CONTAINS BOTH THE EXPANDED DATA AND ORIGINAL SPECTROSCOPY DATA WITH SECOND-ORDER DERIVATIVE ADDING AUXILIARY INFORMATION

Strategy	Strategy I								
LV number	Calibration ($n = 1340$)			Valid ($n = 17$)			test ($n = 17$)		
	r^2	RMSE mg kg^{-1}	RPIQ	r^2	RMSE mg kg^{-1}	RPIQ	r^2	RMSE mg kg^{-1}	RPIQ
2	0.528	3.740	1.604	0.559	3.928	1.782	-0.253	6.152	0.975
5	0.852	2.093	2.867	0.854	2.263	3.093	-0.118	5.810	1.032
8	0.957	1.136	5.284	0.957	1.230	5.691	-0.219	6.066	0.988

TABLE V

EFFECT OF THE LV NUMBER AND THE PRETREATMENT OF STRATEGY II ON THE EVALUATION CRITERIA VALUES. STRATEGY II INDICATES THAT THE TYPES OF INPUT VARIABLES CONTAINS ORIGINAL SPECTROSCOPY WITH SECOND-ORDER DERIVATIVE ADDING AUXILIARY INFORMATION

Strategy	Strategy II								
LV number	Calibration ($n = 134$)			Valid ($n = 17$)			test ($n = 17$)		
	r^2	RMSE mg kg^{-1}	RPIQ	r^2	RMSE mg kg^{-1}	RPIQ	r^2	RMSE mg kg^{-1}	RPIQ
2	0.979	0.806	7.447	0.176	4.902	1.224	-0.186	5.458	1.099
5	0.883	1.896	3.164	0.143	4.999	1.200	-0.209	5.513	1.088
8	0.462	4.078	1.471	-0.094	5.649	1.062	-0.403	5.936	1.011

TABLE VI

EFFECT OF THE LV NUMBER AND THE PRETREATMENT OF STRATEGY III ON THE EVALUATION CRITERIA VALUES. STRATEGY III INDICATES THAT THE TYPES OF INPUT VARIABLES CONTAINS ORIGINAL SPECTROSCOPY WITH SECOND-ORDER DERIVATIVE AND PCA METHOD ADDING AUXILIARY INFORMATION

Strategy	Strategy III								
LV number	Calibration ($n = 134$)			Valid ($n = 17$)			test ($n = 17$)		
	r^2	RMSE mg kg^{-1}	RPIQ	r^2	RMSE mg kg^{-1}	RPIQ	r^2	RMSE mg kg^{-1}	RPIQ
2	0.823	0.002	2.551	0.827	0.002	2.807	0.407	0.004	1.488
5	0.924	0.002	3.882	0.888	0.002	3.491	0.184	0.005	1.269
8	0.963	0.001	5.572	0.945	0.001	4.969	-0.116	0.006	1.085

model performance when Strategy III was utilized as the pre-processing method (Table VI). The model performed best when the LV was equal to 2, with R^2 values of 0.407, RMSE values of 0.004 and RPIQ values of 1.488. The model over-fitted as the number of LVs increased and was most severe at a value of 8 for LVs, when the value of R^2 for the test set was -0.116, the value of RMSE was 0.006 and the value of RPIQ was 1.085.

G. Four CNN Models

The first DA-CNN model (Model 5) is shown in Fig. 11(a). It is a deep learning model with 2 convolutional layers and a fully connected layer. Each convolutional layer is configured with a ReLU activation function, and after the two convolutional layers, there is a flattened layer and a fully connected layer. The second convolutional neural network model (Model 6) is shown in Fig. 11(b). It includes three convolutional layers, and there is a ReLU activation function after the first two convolutional layers. In addition to configuring a ReLU activation function for the third convolutional layer, there is also a flattened layer and a fully connected layer. Compared

with the first layer model, the second layer model adds a layer with a 3*3 convolution layer on the basis of the first two layers of the convolution layer, and the full connection is no longer a 2128-dimensional tensor to a 1-dimensional tensor. The full connection is from 2112 dimensions to a 1 dimension, with an extra convolutional layer, this is equivalent to adding a different size filter to learn the spectral data of the feature in order to get more information. The third convolutional neural network model has three convolutional layers and is composed of two fully connected layers. The third model is shown in Fig. 11(c). It has one more convolutional layer and one fully connected layer than the first model. There is a ReLU function behind each convolutional layer. Table VII shows that the training set, validation set and test set of Model 6 are with the better RPIQ, and the smaller RMSE value. This indicates that the more complex soil hyperspectral dataset and more convolutional layers can extract more one-dimensional soil spectral signals. Different levels of features make the features richer and more complex, thereby generating a better regression model. After Model 7 is deepened in the fully connected layer, although the nonlinear expression ability of the model is improved to increase the learnability of the

TABLE VII

STATISTICAL SUMMARY OF CU ESTIMATION WITH THE PRETREATMENT OF STRATEGY DEVELOPED FROM MODEL 1, 2, 3 AND 4. STRATEGY INDICATES THAT THE TYPES OF INPUT VARIABLES CONTAINS BOTH THE EXPANDED DATA AND ORIGINAL SPECTROSCOPY DATA WITH SECOND-ORDER DERIVATIVE ADDING AUXILIARY INFORMATION

Strategy	Strategy I								
Models	Calibration ($n = 1340$)			Valid ($n = 17$)			test ($n = 17$)		
	r^2	RMSE mg kg^{-1}	RPIQ	r^2	RMSE mg kg^{-1}	RPIQ	r^2	RMSE mg kg^{-1}	RPIQ
Model 5	0.652	3.212	1.868	0.604	3.724	1.879	0.647	3.267	1.836
Model 6	0.762	2.656	2.259	0.752	2.947	2.376	0.761	2.687	2.233
Model 7	0.749	2.723	2.203	0.737	3.036	2.305	0.748	2.756	2.177
Model 8	0.941	1.322	4.537	0.949	1.341	5.222	0.942	1.324	4.531

model, after adding 2 layers, it learns more interference in the soil to correctly predict the spectral information of the Cu concentration, which make the results overfitting. And thus, the model performance is better than that of the first model, but it is worse than that of the second model. Model 8 [see Fig. 11(d)] adds residual connections on the basis of model 7, which improves the learning ability of network representation and predicts heavy metal content best. The above findings further illustrate the pervasiveness of data augmentation for deep learning.

H. Impact of Two Different Traditional Pretreatments on Deep Learning Models

In order to explore the impact of traditional preprocessing methods on deep learning, two different traditional preprocessing methods were used and combined with deep learning models to perform regression prediction on Cu content. As can be seen from Table VIII, the two deep learning models with the best performance in strategy I, under the preprocessing methods of strategy II and strategy III, the accuracy of model prediction is no better than that under the condition of strategy I preprocessing. For the two deep learning models, the effect of using strategy II preprocessing is significantly better than using Strategy III. The RMSE, R^2 and RPIQ of model 8 in strategy II pretreatment test set were 4.316, 0.259 and 1.390, respectively. In Strategy III, RMSE, R^2 and RPIQ were 0.007, -1.099 and 0.826, respectively. In model 4, the values of RMSE, R^2 and RPIQ in the test set of model Strategy II are 4.446, 0.213 and 1.349, respectively. However, in Strategy III, the effect was poor, and R^2 was negative, while RMSE and RPIQ were 0.005 and 1.089, respectively.

V. DISCUSSION

In this paper, we studied the data after processing the measured hyperspectral data combined with second derivative processing and the data enhancement method. In addition to using the measured spectral data, we added the 2D index and other heavy metals and SOM in the soil 1 features using a deep learning model to predict the Cu concentration. After second-order derivation of the spectral data, a more detailed relationship between spectral band information and the Cu concentration can be obtained. Using spectral parameters that are more sensitive to soil Cu, we experimented with two-band

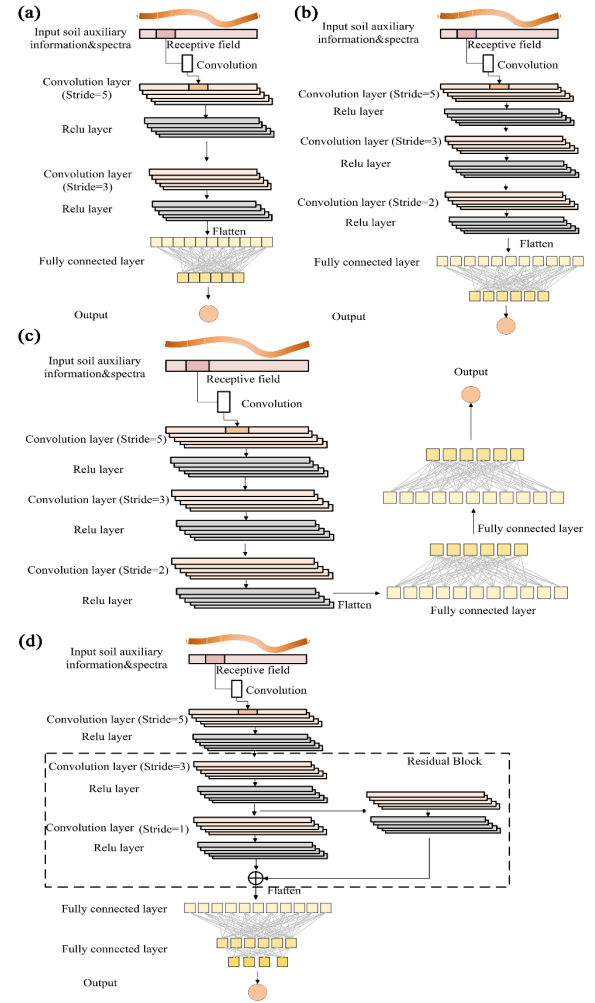


Fig. 11. The structure of the CNN model. (a) Model 5, (b) Model 6, (c) Model 7, (d) Model 8.

spectral indices called RI and NDI to detect their correlation with the soil Cu concentration. Many existing studies have confirmed that certain hyperspectral bands are related to other substances in the soil, but the influence of the composition of this substance on the bands changes with the geographic location; therefore to understand the prediction mechanism of Cu in the Eastern Junggar District, we introduced the concentration of five substances (SOM, As, Cr, Pb, and Zn)

TABLE VIII

STATISTICAL SUMMARY OF CU ESTIMATION WITH THE PRETREATMENT OF STRATEGY II AND STRATEGY III DEVELOPED FROM MODEL 4 AND 8, RESPECTIVELY. STRATEGY II INDICATES THAT THE TYPES OF INPUT VARIABLES CONTAINS ORIGINAL SPECTROSCOPY WITH SECOND-ORDER DERIVATIVE ADDING AUXILIARY INFORMATION. STRATEGY III INDICATES THAT THE TYPES OF INPUT VARIABLES CONTAINS ORIGINAL SPECTROSCOPY WITH SECOND-ORDER DERIVATIVE AND PCA METHOD ADDING AUXILIARY INFORMATION

Strategies	Models	Calibration ($n = 134$)			Valid ($n = 17$)			test ($n = 17$)		
		r^2	RMSE mg kg^{-1}	RPIQ	r^2	RMSE mg kg^{-1}	RPIQ	r^2	RMSE mg kg^{-1}	RPIQ
Strategy II	Model8	0.541	3.768	1.593	0.161	4.948	1.213	0.259	4.316	1.390
Strategy III	Model 8	0.905	0.002	3.498	-0.736	0.007	0.843	-1.099	0.007	0.826
Strategy II	Model 4	0.493	3.959	1.516	0.244	4.695	1.278	0.213	4.446	1.349
Strategy III	Model 4	-0.250	0.006	0.965	-0.264	0.006	0.988	-0.207	0.005	1.089

to assist in predicting the Cu concentration. Finally, we used the extended spectral data processed by the second-order differential, the 8 data selected by the 2D index, and the spectral concentration of the five soil substances as features and used a deep learning model that can learn the fine-grained spectral features to predict the Cu concentration in the soil.

Deep learning has been reasonably well developed within the available spectroscopy domain, and this research chose to employ a deep learning model rather than the more commonly applied PLS model when using hyperspectral empirical data to predict the Cu concentration in the soil. Nevertheless, the quality of feature selection can have a significant impact on the performance of a deep learning model, and the use of well-chosen features as input to a deep neural network enables a more effective model for predicting the Cu concentration in the soil. Consequently, to achieve more effective outcomes in the deep learning model for predicting the Cu concentration in the soil, the present experiment adopted the second-order derivative preprocessing methodology to preprocess the hyperspectral data by learning from the experience of previous scholars on the processing and selection of features, and incorporated two-band spectral indices as well as five substances (SOM, As, Cr, Pb, and Zn) as auxiliary information to predict the Cu concentration in the soil. As can be observed from the outcomes, the features selected for this research demonstrated superior performance in both deep learning models. Furthermore, of the selected auxiliary data in the better performing and interpretable TabNet model, only one ranked outside the top 30 in terms of the importance of the features, which indicates that the auxiliary data constitute the more significant information in the improved TabNet model and that the auxiliary information has a positive impact on the model performance. While 18 of the soil hyperspectral data processed by the second-order derivative algorithm were ranked in the top 30 in terms of the importance of the features, which demonstrates that a number of spectral bands processed by the second-order derivative are essential to the performance of the model. Simultaneously, the soil hyperspectral data processed by the second-order derivative algorithm accounted for the largest proportion of the input features. As can be demonstrated by the ultimate optimally performing deep learning model with $\text{RMSE} = 1.324$, $R^2 = 0.942$, and $\text{RPIQ} = 4.531$, it is appropriate to handle the soil hyperspectral data by employing the second-order derivative algorithm.

To be able to perform properly, deep learning requires not only good features, but also extensive data. The FSL may reduce data-intensive work on data collection when obtaining sufficient examples with supervised information is either difficult or impossible. This research makes use of the data category from the three broad categories of data, model, and algorithm in few-shot learning, and transforms the training dataset to obtain a priori knowledge via augmenting the soil hyperspectral training dataset. The measured data in the research area only has 168 data points and too little of a data volume cannot train a neural network model with better effects and strong robustness. Therefore, this experiment used data enhancement techniques commonly employed in deep learning to expand the data to enable the neural network training to obtain the optimal hyperparameters. The values of hyper-parameters were selected for this research by taking the reference from the guidelines for hyperparameters mentioned in the TabNet. For most datasets, $N_{\text{steps}} \in [3, 10]$ is the best option. Nevertheless, the value of N_{steps} should be larger when more feature quantities need to be learned, in which case the network is too deep can lead to potentially-problematic ill-conditioned matrices. Taking into account the small amount of expanded data used in this experiment, and referring to the setting of N_{steps} for smaller datasets in TabNet, the debugging range of N_{steps} for this experiment is 3-5, and a setting of 3 gives the most effective results with other hyperparameters unchanged. N_d and N_a provide a balance between the performance and complexity of the model, with equal values that are reasonable for most datasets. Simultaneously, the values of N_d and N_a should not be too large, or else overfitting and poor generalization may occur. For this experiment, the values of N_d and N_a are still made equal, and the reference set mentioned by TabNet (i.e. $\{8, 16, 24, 32, 64, 128\}$) has been experimented individually, and the final result is that 8 is the optimum value. The value of γ is positively correlated with the value of N_{steps} . As the value of N_{steps} is relatively small, γ is chosen from $\{1.0, 1.2, 1.3, 1.4, 1.5\}$, and the debugging result is that the model performs most effectively when γ is 1.3. In comparison to a small batchsize, when the batchsize is enlarged, it is faster to process the same amount of data, yet when it is too large, the memory tends to be insufficient and it affects the outcome of training loss and generalization ability. According to the volume of soil datasets and the number of features, as well

as referring to the settings of B and B_v values in TabNet, B is chosen from $\{256, 512, 1024, 2048, 4096, 8192\}$ and B_v is chosen from $\{128, 256, 512\}$, and after debugging, it is concluded that $B = 2048$ and $B_v = 128$ are the most suitable parameters. For the present research, the expanded data together with the auxiliary soil information comprised the training set and yielded relatively positive outcomes in both the CNN and TabNet models, yet the outcomes were less favorable in PLSR models. Data augmentation is a commonly applied technique in deep learning, and when coupled with the ResNet block, both the CNN and TabNet models can effectively learn useful information for the model from a wider range of features, and achieve the optimal outcomes respectively. This illustrates that the use of data augmentation techniques solely, without analyzing and refining models for soil hyperspectral data, would not achieve optimal outcomes. To optimize the outcome of predicting Cu concentration, it is necessary to take into account data augmentation techniques, as well as the strategies of deep learning networks enhanced by relevant expertise. What also has an impact on the accuracy of the TabNet model is the size of the training dataset, as it increases, the training, validation, and test sets also undergo corresponding changes, with R^2 value, RPIQ value becoming larger and RMSE value becoming smaller. This demonstrates that the increased dataset would have a beneficial effect on the improved TabNet model to learn vital patterns from manifold data sources. When the number or type of hyperspectral data decreases, the deep learning model cannot be fully trained and the optimal results cannot be obtained (Table VIII). And in Table IV, the combination of data augmentation and PLSR performed the worst in the test set, with R^2 even showing a negative number. Yet, when contrasted with the test set, the performance of PLSR in the training and validation sets performed well with the increasing number of the LV, which was only slightly inferior to the optimal model of deep learning, while the model in the test set became worse with the increasing number of LV, demonstrating that the increase in latent variables enabled PLSR to learn more noise information from the expanded data, making the model appear to be overfitted. Both the PCA processing and PLSR combination of hyperspectral data have yielded a high level of heavy metal prediction in previous studies without the use of data augmentation (Table VI). The current experiment also preprocessed the hyperspectral data with a combination of PCA and data augmentation techniques and predicted the Cu concentration by employing PLSR and revealed that the performance was better than the model without PCA processing. The reason for this is that the method of expanding the data is to add random noise, which is beneficial to the training of the deep learning model, yet enables PLSR to learn more noise information. The PCA method can minimize the input of noise information, and the model performs optimally when the number of LV is minimal, since the more LV there are, the more noise in the variables is learned, resulting in poor model performance and a state of overfitting.

Apart from the influence of the preprocessing approach on the ultimate outcomes, the choice and optimization of the model are also vital for the prediction of the Cu concentration

in the soil. The preprocessing approach employed in this research performed at its best in deep learning. In particular, the effects of the CNN model and TabNet model with ResNet learning were excellent; whereas the CNN model with ResNet learning was slightly higher than the modified TabNet model in terms of its performance, yet its interpretability was weaker, while the TabNet model can directly select the more significant features and rank them in order of importance. In previous research on the prediction of the heavy metal concentration from soil hyperspectral, Pyo et al. [27] performed PCA on 98 visible and near-infrared spectroscopy soil samples and experimentally revealed that the CNN model with convolutional autoencoder yielded the highest As, Cu, and Pb estimates. For this research, the PCA was used to reduce the dimensionality and extract the essential information as the data volume was small, which reduced the burden of the neural network to learn excessive features from the original data, and both the CNN model with convolutional autoencoder and the CNN model with ResNet learning added proposed in this paper showed satisfactory results in predicting the heavy metal concentration. This indicates that CNN and its adapted models offer great potential in predicting the heavy metal concentration with high accuracy, simplicity, and rapid prediction by utilizing soil hyperspectral data. Nonetheless, none of the above models are able to explain the impact of the input data on the prediction of heavy metal concentration well enough nor are they able to carry out the analysis related to soil spectroscopy due to the limitations of convolutional neural networks, which, however, still have a good rapid detection capability. The deep learning model mentioned in this paper is mainly used to extract important information and to predict the Cu concentration. There are, nevertheless, existing studies that can apply deep learning models to data augmentation and perform regression predictions with traditional regression models, which are distinct from the purpose of the deep learning models used in this paper [28]. Also, some studies were done without the use of real soil hyperspectral data and instead used some data from other mediums combined with deep learning to estimate heavy metal concentration. Since no hyperspectral data were used and the number of features used was relatively small, yet with a strong correlation to the heavy metal concentration, little but essential information can contribute to a better performance of the deep learning model GRU [29]. This is consistent with the improved TabNet model, whereby the auxiliary information is all ranked in the top 30 of the importance list. Compared to the sensitivity analysis where the SHAP values represent the contribution of the input variables to the predicted concentrations of the four heavy metals, TabNet only lists 30 significant features in the model, however, the model is able to obtain the magnitude of the contribution of all the features to the prediction of the heavy metals. The TabNet model has not only a regression function, but also a sensitivity analysis function, which is both comprehensive and powerful. The improved TabNet model with its multiple functions offers both high accuracy and high speed compared to deep learning models that require large-scale data training [31], and is suitable for practical application in the Zhundong-Xinjiang Economic & Technological Development Zone.

TABLE IX
SOIL ABSORBANCE AND AUXILIARY INFORMATION IN THE VISIBLE-NEAR-INFRARED REGIONS

Ranking	Soil constituent or Wavelength	Relative importance of the predictor variable	Reference
1	Cr	0.05932	Balasoju et al. (2001) [68]
2	SDR ₆₈₈ , SDR ₁₈₅₁	0.05888	Hong et al. (2018) [42]
3	As	0.05646	Nearing et al. (2014) [69], Kabirinejad et al. [70]
4	SOM	0.05380	Wang et al. (2018) [71], Sharma et al. (2015) [72], Zhang et al. (2018) [73]
5	897 nm	0.05185	Sherman and Waite (1985) [61], Wang et al. (2011) [62]
6	SDR ₄₀₃ , SDR ₂₁₉₅	0.05002	Hong et al. (2018) [42]
7	1909 nm	0.04617	Liu et al. (2018) [67], Hunt (1977) [57]
8	2036 nm	0.03288	Clark et al. (1990) [59], Clark (1999) [58]
9	2152 nm	0.03058	Clark et al. (1990) [59], Clark (1999) [58]
10	1761 nm	0.03045	Ben-Dor et al. (1997) [55]
11	631 nm	0.02843	Ben-Dor et al. (1997) [55]
12	SDR ₁₅₈₁ , SDR ₂₃₃₁	0.02820	Hong et al. (2018) [42]
13	Zn	0.02726	Li (2007) [74], Mittal et al. (2015) [75], Abd-Elfattah et al. (1981) [76]
14	Pb	0.02722	Li (2007) [74], Mittal et al. (2015) [75], Abd-Elfattah et al. (1981) [76]
15	2365 nm	0.02396	Clark (1999) [58]
16	SDR ₆₂₃ , SDR ₂₀₉₁	0.02291	Hong et al. (2018) [42]
17	2380 nm	0.02196	Fourty et al. (1996) [56]
18	SDR ₈₅₁ , SDR ₈₈₀	0.02090	Hong et al. (2018) [42]
19	2284 nm	0.02022	Clark et al. (1990) [59], Ben-Dor et al. (1997) [55] Post and Noble (1993) [60]
20	1829 nm	0.01728	Hunt (1977) [57]
21	SDR ₁₉₄₉ , SDR ₂₁₉₅	0.01649	Hong et al. (2018) [42]
22	2320 nm	0.01637	Clark (1999) [58]
23	2266 nm	0.01627	Clark et al. (1990) [59]
24	2327 nm	0.01612	Ben-Dor et al. (1997) [55]
25	1947 nm	0.01607	Fourty et al. (1996) [56]
26	SDR ₈₇₀ , SDR ₂₃₈₅	0.01467	Hong et al. (2018) [42]
27	2321 nm	0.01430	Clark (1999) [58]
28	1985 nm	0.01214	White (1971) [66], Hunt (1977) [57]
29	972 nm	0.01188	Liu et al. (2013) [63], Scheinost et al. (1998) [64]
30	1422 nm	0.01122	Clark et al. (1990) [59], Oinuma and Hayashi (1965) [65]

The improved TabNet model achieves the best results through expanded data and additional auxiliary information. According to the advantages of the TabNet, the top 30 most important features input to the neural network are sorted, as shown in Table IX. It can be seen from Table IX that except for some bands selected from the hyperspectral bands, only one of the auxiliary information selected in this experiment is not in the table, indicating that this information is an important feature and is useful for the generation of deep learning models.

In the ranking of importance, there are two heavy metal elements in the top 5 features. This is because when two or more heavy metals coexist, the heavy metals formed by physicochemical reactions between each other and between metals and soil ions will affect the physicochemical environment of the soil, and hence will be reflected in the spectrum measurement. Since peat carbon has a greater correlation with the retention of Cu and Cr in the soil and soil organic matter has similar effects on the concentration of Cr and Cu, the correlation between Cu and Cr is the largest in the Eastern Junggar coalfield [68]. Copper has the following five forms

in soil, with increased levels of exogenous heavy metals, the various forms of heavy metals are redistributed in the soil because Cu and As are in the soil of the same environment, and the main form of Cu in the soil is similar to that of As; therefore, As and Cu have strong relevance [69], [70]. Copper ions have a special affinity for organic matter, because Cu^{2+} easily forms complexes, so Cu has a strong correlation with organic matter [71], [72], [73]. Both Zn and Pb are shown to correlate with Cu, and studies have demonstrated that as the concentration of exogenous Pb^{2+} , Cu^{+} , Zn^{2+} , and Cd^{2+} increases, the chance of their collision with the soil surface increases correspondingly, thus improving soil sorption of heavy metals [74], [75], [76].

Seven of the eight 2D index data selected are included in the top 30 features of importance, proving that the data pairs selected by Pearsons coefficient are also extremely valuable in the overall characteristics [42]. In this experiment, the order of importance of 2D data in the features is not exactly the same as the order of Pearsons coefficient, suggesting that although features of higher importance can be identified by the magnitude of Pearsons coefficient, new weights are judged

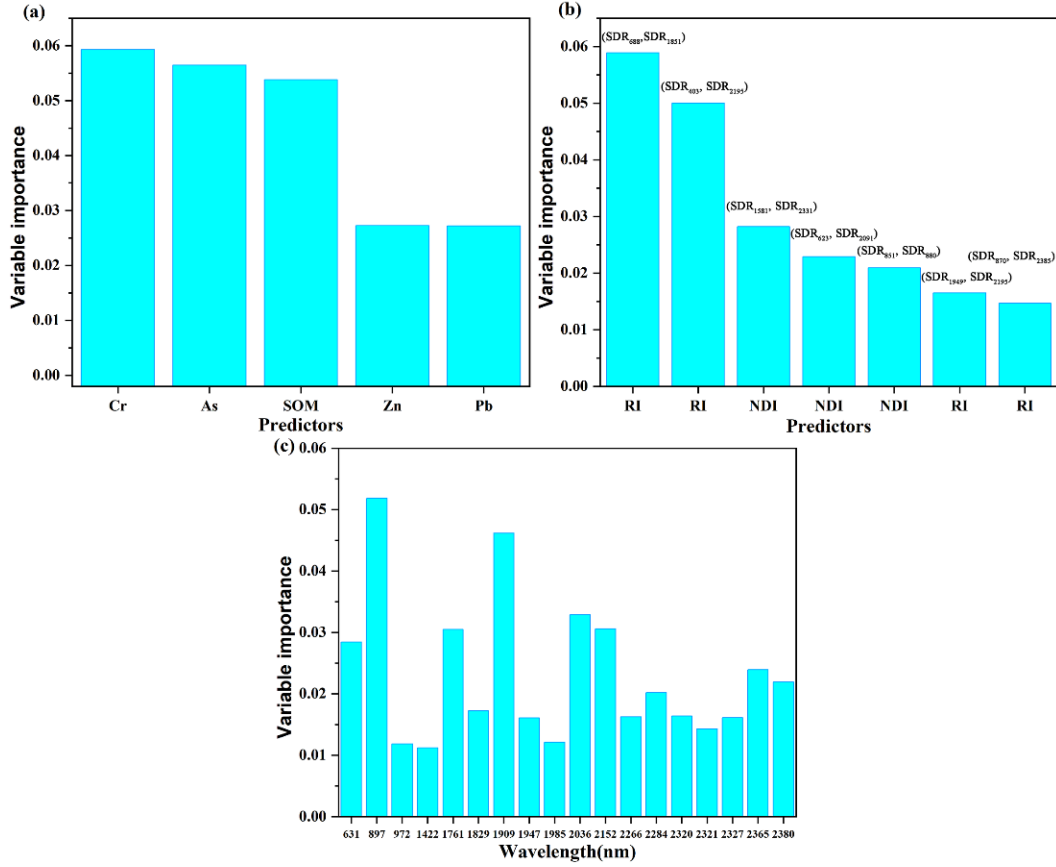
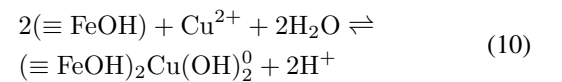
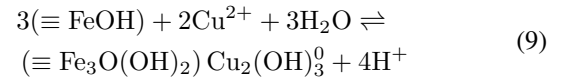


Fig. 12. Relative importance of the predictor variables used in model 4.

when the model uses the data to find the optimal combination of features, and will no longer be based on a single index. In this experiment, the performance of the ratio index is better than that of the normalized difference index. The four sets of data selected by the ratio index are all in the top 30 features, with 1451 nm ranking 2nd and 3 nm ranking 6th, and the normalized difference index has only three features among the top 30 feature selections; the highest ranking of the data selected using the normalized difference index is 12.

The 18 bands processed by the second derivative are among the top 30 features selected, proving that the spectral features after second-derivative processing have a greater influence on the deep learning model. In the top 10 bands, because 897 nm is at the Fe-OH [61] stretching vibration peak, with increasing Cu^+ concentration, an Fe-O-(Cu) structure is formed on the surface of hematite, and the adsorption mechanism is approximately as shown in formulas 3 and 4 [62]. The significant correlation at 1909 nm is the result of the hydroxyl group and the Al-OH lattice structure [57], [67]. The bands at 2036 nm, 2152 nm, and 1761 nm are due to organic molecules and proteins with C=O, Al-OH, and C-H bonds in soil organic carbon [55], [58], [59], [65]. The bands ranked 11th to 20th corresponding to the substance that are all animal and plant residues and intermediate products; this is due to the reaction of Cu^{2+} with organic functional groups, which fixes the metal ion in a stable inner complex [64]. There are a total of 8 bands with an importance ranking between 21

and 30. 2320 nm [58] and 2321 nm [58] are the corresponding bands of methyl groups where the soil components containing methyl groups can promote the adsorption of Cu^+ , and the material corresponding to the bands of 2327 nm [55] and 1947 nm [56] contains some coordination groups, which can fix copper by complexation. The material corresponding to the 2266 nm band [59] is gibbsite, which produces a large amount of H^+ when hydrolyzed, affecting the amount of Cu^+ [77]. The strong response of Cu corresponding to the carbonate component of the 1985 nm band can be attributed to the fact that Cu lies in the carbonate-bound state. Because Cu^{2+} can form strong complexes on the surface when it reacts with ferrihydrite, it has a strong correlation at 972 nm [63], [64]. The relatively high degree of water in the 1422 nm band is due to the combined effect of water molecular vibrations and OH^- [59], [65].



In this experiment, in the data processing stage, although data enhancement technology is used, there are still fewer methods for soil hyperspectral data enhancement. Compared with image field data enhancement technology, it is not mature and complete. Despite the generalization performance of small

datasets through a combination of deep learning and few-shot learning, the use of few-shot learning methods and the use of empirical hyperspectral data to predict heavy metal concentration is unprecedented in previous studies and requires data augmentation with additional knowledge of soil science and deep learning. And, existing studies that use empirical hyperspectral data to predict heavy metal concentration, while providing high accuracy, are unable to provide a reasonable interpretation of the input features. More studies with a combination of few-shot learning and deep learning are needed in empirical soil hyperspectral data, as such studies not only provide higher accuracy, but also the interpretation of the features used. In addition, although the improved CNN model and TabNet model work well, other deep learning models used for comparison have poor results. The focus of future research is to develop a model that is more suitable for soil hyperspectra and to try to avoid deep learning models that are not suitable for specific soil hyperspectra.

VI. CONCLUSION

In this study, we can prove that the combination of soil near-infrared spectroscopy and auxiliary information in soil can be exploited to predict the copper concentration using deep learning methods. Due to the low concentration of Cu in the soil, this study not only uses the spectral information processed by the second derivative but also expands the spectral data and adds soil auxiliary information to increase the abundance of the effective information. In sharp contrast, when predicting the Cu concentration, only a very small number of hyperspectral bands are ranked in the top 30 in terms of importance. Experiments have proven that organic matter, Fe oxides, and clay minerals significantly affect the response of near-infrared spectroscopy to the Cu concentration, and the addition of features such as the 2D index and other heavy metal concentrations is also important. The abovementioned characteristics provide effective information for the deep learning model to train a model with high accuracy and strong robustness. The improved TabNet and CNN high-precision regression prediction results show that using near-infrared hyperspectral imaging and auxiliary information in soil to identify heavy metal pollution is an effective method.

Using hyperspectral data for Cu concentration detection can detect pollution in a timely manner, avoiding further deterioration of heavy metal pollution and enabling soil remediation measures to be taken. Due to the different compositions of soil in different regions, different spectral characteristics are needed to model in heavy metal pollution areas, and then, effective pollution prevention measures can be taken. In recent years, due to the wide application of ultralight and high-resolution drones and the successive launches of hyperspectral satellites at home and abroad, this research can facilitate large-scale monitoring and prevention measures for heavy metal-polluted areas.

REFERENCES

[1] A. McBratney, D. J. Field, and A. Koch, "The dimensions of soil security," *Geoderma*, vol. 213, pp. 203–213, 2014.

[2] R. Zhang, S. Liu, and S. Zheng, "Characterization of nano-to-micron sized respirable coal dust: Particle surface alteration and the health impact," *Journal of Hazardous Materials*, vol. 413, p. 125447, 2021.

[3] A. Zj, B. Jwa, A. Rw, and W. A. Ping, "Using multi-fractal analysis to characterize the variability of soil physical properties in subsided land in coal-mined area," *Geoderma*, vol. 361, p. 114054, 2020.

[4] L. Song, J. Jian, D. J. Tan, H. B. Xie, Z. F. Luo, and B. Gao, "Estimate of heavy metals in soil and streams using combined geochemistry and field spectroscopy in wan-sheng mining area, chongqing, china," *International Journal of Applied Earth Observation and Geoinformation*, vol. 34, pp. 1–9, 2015.

[5] Y. Wu, J. Chen, J. Ji, G. Peng, Q. Liao, Q. Tian, and H. Ma, "A mechanism study of reflectance spectroscopy for investigating heavy metals in soils," *Soil Science Society of America Journal*, vol. 71, no. 3, pp. 918–926, 2007.

[6] R. Hong-Yan, D. Zhuang, A. N. Singh, P. Jian-Jun, Q. Dong-Sheng, and S. Run-He, "Estimation of as and cu contamination in agricultural soils around a mining area by reflectance spectroscopy: a case study," *Pedosphere*, vol. 19, no. 6, pp. 719–726, 2009.

[7] Y. Song, F. Li, Z. Yang, G. A. Ayoko, R. L. Frost, and J. Ji, "Diffuse reflectance spectroscopy for monitoring potentially toxic elements in the agricultural soils of changjiang river delta, china," *Applied Clay Science*, vol. 64, pp. 75–83, 2012.

[8] L. Cai, Z. Xu, M. Ren, Q. Guo, X. Hu, G. Hu, H. Wan, and P. Peng, "Source identification of eight hazardous heavy metals in agricultural soils of huizhou, guangdong province, china," *Ecotoxicology & Environmental Safety*, vol. 78, pp. 2–8, 2012.

[9] W. Sun, X. Zhang, X. Sun, Y. Sun, and Y. Cen, "Predicting nickel concentration in soil using reflectance spectroscopy associated with organic matter and clay minerals," *Geoderma*, vol. 327, pp. 25–35, 2018.

[10] Y. Hong, R. Shen, H. Cheng, S. Chen, and Y. Liu, "Cadmium concentration estimation in peri-urban agricultural soils: Using reflectance spectroscopy, soil auxiliary information, or a combination of both?" *Geoderma*, vol. 354, p. 113875, 2019.

[11] C. Hang, R. Shen, Y. Chen, Q. Wan, T. Shi, J. Wang, W. Yuan, Y. Hong, and X. Li, "Estimating heavy metal concentrations in suburban soils with reflectance spectroscopy," *Geoderma*, vol. 336, pp. 59–67, 2019.

[12] M. S. Luce, N. Ziadi, B. Gagnon, and A. Karam, "Visible near infrared reflectance spectroscopy prediction of soil heavy metal concentrations in paper mill biosolid- and liming by-product-amended agricultural soils," *Geoderma*, vol. 288, pp. 23–36, 2017.

[13] T. Pfister, J. Charles, and A. Zisserman, Eds., *Domain-adaptive discriminative one-shot learning of gestures*. Cham: Springer International Publishing, 2014.

[14] O. Vinyals, C. Blundell, T. Lillicrap, and D. Wierstra, "Matching networks for one shot learning," *Advances in Neural Information Processing Systems*, vol. 29, pp. 3630–3638, 2016.

[15] Y. Wang, Q. Yao, J. T. Kwok, and L. M. Ni, "Generalizing from a few examples," *ACM Computing Surveys*, vol. 53, no. 3, pp. 1–34, 2020.

[16] H. Qi, M. Brown, and D. G. Lowe, "Low-shot learning with imprinted weights," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Salt Lake City, UT, USA, 18–23 June 2018, pp. 5822–5830.

[17] B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum, "Human-level concept learning through probabilistic program induction," *Science*, vol. 350, no. 6266, pp. 1332–1338, 2015.

[18] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap, "Meta-learning with memory-augmented neural networks," in *Proceedings of the 33rd International Conference on Machine Learning*, New York, NY, USA, 2016.

[19] E. J. Bjerrum, M. Ghlader, and T. Skov, "Data augmentation of spectral data for convolutional neural network (cnn) based deep chemometrics," *arXiv preprint arXiv:1710.01927*, 2017.

[20] J. Wang, S. Kim, and Y. Lee, "Speech augmentation using wavenet in speech recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, Brighton, UK, 12–17 May 2019, p. 18777935.

[21] M. Wu, S. Wang, S. Pan, A. C. Terentis, and X. Zhu, "Deep learning data augmentation for raman spectroscopy cancer tissue classification," *Scientific Reports*, vol. 11, no. 8, p. 23842, 2021.

[22] J. M. Haut, M. E. Paoletti, J. Plaza, A. Plaza, and J. Li, "Hyperspectral image classification using random occlusion data augmentation," *IEEE Geoscience and Remote Sensing Letters*, vol. 15, no. 11, pp. 1751–1755, 2019.

[23] X. L. Zhang, T. Lin, J. F. Xu, X. Luo, and Y. B. Ying, "Deepspectra: an end-to-end deep learning approach for quantitative spectral analysis," *Analytica Chimica Acta*, vol. 1058, pp. 48–57, 2019.

- [24] X. Yu, H. Lu, and D. Wu, "Development of deep learning method for predicting firmness and soluble solid content of postharvest korla fragrant pear using vis/nir hyperspectral reflectance imaging," *Postharvest Biology and Technology*, vol. 141, pp. 39–49, 2018.
- [25] X. Yu, J. Wang, S. Wen, J. Yang, and F. Zhang, "A deep learning based feature extraction method on hyperspectral images for nondestructive prediction of tvb-n content in pacific white shrimp (*litopenaeus vannamei*)," *Biosystems Engineering*, vol. 178, pp. 244–255, 2019.
- [26] J. Liu, M. Osadchy, L. Ashton, M. Foster, C. J. Solomon, and S. J. Gibson, "Deep convolutional neural networks for raman spectrum recognition: a unified solution," *Analyst*, vol. 142, no. 21, pp. 4067–4074, 2017.
- [27] J. C. Pyo, S. M. Hong, Y. S. Kwon, M. S. Kim, and K. H. Cho, "Estimation of heavy metals using deep neural network with visible and infrared spectroscopy of soil," *Science of the Total Environment*, vol. 74, no. 1, p. 140162, 2020.
- [28] X. Zhou, S. Jun, T. Yan, L. Bing, Y. Y. Hang, and Q. S. Chen, "Hyperspectral technique combined with deep learning algorithm for detection of compound heavy metals in lettuce," *Food Chemistry*, vol. 321, no. 8, p. 126503, 2020.
- [29] Y. Q. Jiang, C. L. Li, H. X. Song, and W. H. Wang, "Deep learning model based on urban multi-source data for predicting heavy metals (cu, zn, ni, cr) in industrial sewer networks," *Journal of Hazardous Materials*, vol. 432, no. 2, p. 128732, 2022.
- [30] H. L. Qiao, X. B. Shi, H. Z. Chen, J. Y. Lyu, and S. Y. Hong, "Effective prediction of soil organic matter by deep svd concatenation using ft-nir spectroscopy," *Soil & Tillage Research*, vol. 215, no. 9, p. 105223, 2022.
- [31] L. Zhao, X. Guo, Z. Xu, and M. Ding, "Soil properties: Their prediction and feature extraction from the lucas spectral library using deep convolutional neural networks," *Geoderma*, vol. 402, no. 7, p. 115366, 2021.
- [32] S. O. Arik and T. Pfister, "Tabnet: attentive interpretable tabular learning," *arXiv*, 2020.
- [33] A. Abliz, Q. Shi, M. Keyimu, and R. Sawut, "Spatial distribution, source, and risk assessment of soil toxic metals in the coal-mining region of northwestern china," *Arabian Journal of Geosciences*, vol. 11, no. 24, p. 793, 2018.
- [34] J. Xia, "The study of remote sensing dynamic monitoring for coalfield fire area in shuixigou, xinjiang," *IOP conference series. Earth and environmental science*, vol. 17, no. 1, p. 12097, 2014.
- [35] A. Savitzky and M. J. Golay, "Smoothing and differentiation of data by simplified least squares procedures," *Analytical Chemistry*, vol. 36, no. 8, pp. 1627–1639, 1964.
- [36] Z. P. Zhang, J. L. Ding, J. Z. Wang, and X. Y. Ge, "Prediction of soil organic matter in northwestern china using fractional-order derivative spectroscopy and modified normalized difference indices," *Catena*, vol. 185, p. 104257, 2020.
- [37] B. Nadia, M. C. Artur, C. Brito, and F. M. Delfim, "A fractional calculus on arbitrary time scales: Fractional differentiation and fractional integration," *Signal Processing*, vol. 107, no. 5, pp. 230–237, 2015.
- [38] P. J. Tong, Y. P. Du, K. Y. Zheng, T. Wu, and J. J. Wang, "Improvement of nir model by fractional order savitzkygolay derivation (fosgd) coupled with wavelength selection," *Chemometrics and Intelligent Laboratory Systems*, vol. 143, no. 2, pp. 40–48, 2015.
- [39] J. Fainberg, P. Bell, M. Lincoln, and S. Renals, "Improving childrens speech recognition through out-of-domain data augmentation," in *Conference: Interspeech*, 2016, pp. 1598–1602.
- [40] Q. Wen, L. Sun, F. Yang, X. Song, J. Gao, X. Wang, and H. Xu, "Time series data augmentation for deep learning: A survey," in *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence (IJCAI-21) Survey Track*, 2021, pp. 4653–4660.
- [41] M. Wu, S. W. Wang, S. R. Pan, A. C. Terentis, J. Strasswimmer, and X. Q. Zhu, "Deep learning data augmentation for raman spectroscopy cancer tissue classification," *Scientific Reports*, vol. 11, no. 8, p. 23842, 2021.
- [42] Y. Hong, S. Chen, Y. Zhang, Y. Chen, L. Yu, Y. Liu, Y. Liu, H. Cheng, and Y. Liu, "Rapid identification of soil organic matter level via visible and near-infrared spectroscopy: Effects of two-dimensional correlation coefficient and extreme learning machine," *Science of The Total Environment*, vol. 644, pp. 1232–1243, 2018.
- [43] M. Knadel, R. A. V. Rossel, and F. Deng, "Visible-near infrared spectra as a proxy for topsoil texture and glacial boundaries," *Soil Science Society of America Journal*, vol. 77, no. 5, pp. 568–579, 2012.
- [44] J. A. M. Demattê, A. C. Dotto, A. F. S. Paiva, M. V. Sato, R. S. D. Dalmolin, and M. D. S. B. d. et al., "The brazilian soil spectral library (bssl): a general view, application and challenges," *Geoderma*, vol. 354, p. 113793, 2019.
- [45] S. Chakraborty, B. Li, S. Deb, S. Paul, D. C. Weindorf, and B. S. Das, "Predicting soil arsenic pools by visible near infrared diffuse reflectance spectroscopy," *Geoderma*, vol. 288, pp. 23–36, 2017.
- [46] V. Khosravi, F. D. Ardejani, S. Yousefi, and A. Aryafar1, "Monitoring soil lead and zinc contents via combination of spectroscopy with extreme learning machine and other data mining methods," *Geoderma*, vol. 318, pp. 29–41, 2018.
- [47] T. Mehmood, K. H. Liland, L. Snipen, and S. Sæbø, "A review of variable selection methods in partial least squares regression," *Chemometrics & Intelligent Laboratory Systems*, vol. 118, pp. 62–69, 2012.
- [48] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *European conference on computer vision*, 2014, pp. 818–833.
- [49] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, 27–30 June 2016, pp. 770–778.
- [50] B. M. Véronique, F. A. Elvira, P. Bernard, R. Jean-Michel, and M. Alex, "Critical review of chemometric indicators commonly used for assessing the quality of the prediction of soil attributes by nir spectroscopy," *Trends in Analytical Chemistry*, vol. 29, no. 9, pp. 1073–1081, 2010.
- [51] L. P. Wilding, "Spatial variability: its documentation, accommodation and implication to soil surveys," in *Soil spatial variability*, Las Vegas NV, Netherlands, 30 November–1 December 1984, pp. 166–194.
- [52] E. R. Stoner, "Physicochemical, site and bi-directional reflectance factor characteristics of uniformly moist soils," M. Eng. thesis, Purdue University, West Lafayette, Indiana, 1979.
- [53] Z. Shi, E. Peltier, and D. L. Sparks, "Kinetics of ni sorption in soils: Roles of soil organic matter and ni precipitation," *Environmental Science & Technology*, vol. 46, no. 4, p. 2212, 2012.
- [54] R. Thilo and R. Jorg, "Modelling the potential mobility of cd, cu, ni, pb and zn in mollic fluvisols," *Environ Geochem Health*, vol. 39, no. 6, pp. 1–14, 2017.
- [55] E. Ben-Dor, Y. Inbar, and Y. Chen, "The reflectance spectra of organic matter in the visible near-infrared and short wave infrared region (400–2500 nm) during a controlled decomposition process," *Remote Sensing of Environment*, vol. 61, no. 1, pp. 1–15, 1997.
- [56] T. Fourty, F. Baret, S. Jacquemoud, G. Schmuck, and J. Verdebout, "Leaf optical properties with explicit description of its biochemical composition: Direct and inverse problems," *Remote Sensing of Environment*, vol. 56, no. 2, pp. 104–117, 1996.
- [57] G. R. Hunt, "Spectral signatures of particulate minerals in the visible and near infrared," *Geophysics*, vol. 42, no. 3, pp. 501–513, 1977.
- [58] R. N. Clark and A. N. Rencz, "Spectroscopy of rocks and minerals, and principles of spectroscopy," *Manual of remote sensing*, vol. 3, pp. 3–58, 1999.
- [59] R. N. Clark, T. King, M. Klejwa, G. A. Swayze, and N. Vergo, "High spectral resolution reflectance spectroscopy of minerals," *Journal of Geophysical Research Solid Earth*, vol. 95, no. B8, pp. 12653–12680, 1990.
- [60] J. L. Post and P. N. Noble, "The near-infrared combination band frequencies of dioctahedral smectites, micas, and illites," *Clays and Clay Minerals*, vol. 41, no. 6, pp. 639–644, 1993.
- [61] D. M. Sherman and T. D. Waite, "Electronic spectra of fe³⁺ oxides and oxide hydroxides in the near ir to near uv," *American Mineralogist*, vol. 70, no. 11, pp. 1262–1269, 1985.
- [62] C. L. L. J. J. Z. S. Wang, N. Wang and S. Dou, "Ftir spectroscopic analysis of cu²⁺ adsorption on hematite and bayerite," *Spectroscopy & Spectral Analysis*, vol. 31, no. 9, pp. 2403–2406, 2011.
- [63] B. Liu, D. Qu, X. Chen, Q. Li, and L. Peng, "Effects of flooding and ferrihydrite on copper fractionation in paddy soil," *Procedia Environmental Sciences*, vol. 18, pp. 135–142, 2013.
- [64] A. C. Scheinost, V. B. A. Chavernas, and J. Torrent, "Use and limitations of second-derivative diffuse reflectance spectroscopy in the visible to near-infrared range to identify and quantify fe oxide minerals in soils," *Clays & Clay Minerals*, vol. 46, no. 5, pp. 528–536, 1998.
- [65] K. Oinuma and H. Hayashi, "Infrared study of mixed-layer clay minerals," *American Mineralogist*, vol. 50, no. 9, pp. 1213–1227, 1965.
- [66] W. B. White, "Infrared characterization of water and hydroxyl ion in the basic magnesium carbonate minerals," *American Mineralogist*, vol. 56, no. 1–2, pp. 46–53, 1971.
- [67] Y. Liu, Z. Shi, G. Zhang, Y. Chen, S. Li, Y. Hong, T. Shi, J. Wang, and Y. Liu, "Application of spectrally derived soil type as ancillary data to improve the estimation of soil organic carbon by using the chinese soil vis-nir spectral library," *Remote Sensing*, vol. 10, no. 11, p. 1747, 2018.

- [68] C. F. Balasoiu, G. J. Zagury, and L. Deschênes, "Partitioning and speciation of chromium, copper, and arsenic in cca-contaminated soils: influence of soil composition," *The Science of the total environment*, vol. 280, no. 1, pp. 239–255, 2001.
- [69] M. M. Nearing, I. Koch, and K. J. Reimer, "Complementary arsenic speciation methods: a review," *Spectrochimica Acta Part B: Atomic Spectroscopy*, vol. 99, pp. 150–162, 2014.
- [70] S. Kabirinejad, M. Kalbasi, A. H. Khoshgoftarmansh, M. Hoodaji, and M. Afyuni, "Chemical forms and phytoavailability of copper in soil as affected by crop residues incorporation," *American Journal of Analytical Chemistry*, vol. 5, no. 9, pp. 604–612, 2014.
- [71] F. Wang, J. Yu, Z. Zhang, Y. Xu, and R. A. Chi, "An amino-functionalized ramie stalk-based adsorbent for highly effective Cu^{2+} removal from water: Adsorption performance and mechanism," *Process Safety & Environmental Protection*, vol. 117, pp. 511–522, 2018.
- [72] B. D. Sharma, J. S. Brar, J. K. Chanay, P. Sharma, and P. K. Singh, "Distribution of forms of copper and their association with soil properties and uptake in major soil orders in semi-arid soils of punjab, india," *Communications in Soil Science & Plant Analysis*, vol. 46, no. 4, pp. 511–527, 2015.
- [73] X. Zhang, J. Li, D. Wei, B. Li, and Y. Ma, "The solid-solution distribution of copper added to soils: influencing factors and models," *Journal of Soils and Sediments*, vol. 18, no. 9, pp. 1–10, 2018.
- [74] Y. P. Li, "Adsorption and desorption characteristics of lead, copper, zinc and cadmium in xuzhou and suzhou soils," M. Eng. thesis, Capital Normal University, Beijing, 2007.
- [75] H. Mittal, A. Maity, and S. S. Ray, "The adsorption of Pb^{2+} and Cu^{2+} onto gum ghatti-grafted poly (acrylamide-co-acrylonitrile) biodegradable hydrogel: Isotherms and kinetic models," *The Journal of Physical Chemistry B*, vol. 119, no. 5, pp. 2026–2039, 2015.
- [76] A. Abd-Elfattah and K. Wada, "Adsorption of lead, copper, zinc, cobalt, and cadmium by soils that differ in cation-exchange materials," *European Journal of Soil Science*, vol. 32, no. 2, pp. 271–283, 1981.
- [77] W. T. Ling, "Intercation between charge character and Cu^{2+} adsorption-desorption of soils with permanent variable charge," M. Eng. thesis, Huazhong Agricultural University, Wuhan, 2001.



Yuan Wang received the M.S. degree from the College of Information Science and Engineering, Xinjiang University, Urumqi, China, in 2019. She is currently pursuing the Ph.D. degree with the College of Information Science and Engineering, Xinjiang University, Urumqi, China.

From 2020 to 2022, she is a Visiting Ph.D. Student in Remote Sensing Laboratory, Department of Electronic Engineering, Tsinghua University, Beijing, China, and at the Beijing National Research Center for Information Tsinghua University, Beijing,

China. At the same time, she is a research assistant in Remote Sensing Laboratory, Department of Electronic Engineering, Tsinghua University. Her research interests include remote sensing data processing and machine learning.



Abdugheni Abliz received the B.S. and M.S. degrees from the School of Water resources and Hydropower Engineering, Wuhan University, Wuhan, China, in 2009 and 2012, and the Ph.D. degree from the school of resources and environmental sciences, Xinjiang University, Urumqi, China, in 2016. He is currently an Associate Processor with the School of resources and environmental sciences, Xinjiang University, Urumqi, China. He has authored more than 20 peer-reviewed articles in international journals from multiple domains such as heavy metal pollution, and

remote sensing application. His research interests include soil and plant heavy metal pollution, remote sensing, and their applications in soil science.

Dr. Abdughen has been frequently serving as a reviewer for more than four international journals including the STOTEN, IEEE, ENVIRONMENTAL SCIENCES AND POLLUTION RESEARCH, OPEN GEOSCIENCES.



Hongbing Ma received the B.S. degree in mathematics from Hebei Normal University, Shijiazhuang, China, in 1985, and the M.S. degree in signal and information and the Ph.D. degree in remote sensing from Peking University, Beijing, China, in 1996 and 1999, respectively.

He is currently a Professor with the Department of Electronic Engineering, Tsinghua University, Beijing. His research interests include image processing, pattern recognition, and spatial information processing and its application.



Li Liu received the BSc degree in communication engineering, the MSc degree in photogrammetry and remote sensing, and the PhD degree in information and communication engineering from the National University of Defense Technology (NUDT), China, in 2003, 2005, and 2012, respectively.

She joined the faculty at NUDT in 2012, where she is currently an associate professor with the College of System Engineering. She was a co-chair of seven International Workshops at CVPR, ICCV, and ECCV. She is going to lecture a tutorial at CVPR19.

She was a guest editor of special issues for the IEEE Transactions on Pattern Analysis and Machine Intelligence and the International Journal of Computer Vision. Her current research interests include facial behavior analysis, texture analysis, image classification, object detection, and recognition. Her papers have currently more than 1800 citations in Google Scholar. She currently serves as associate editor of the Visual Computer Journal.



Alishir Kurban received the Ph.D. degree from Xinjiang Institute of Ecology and Geography, Chinese Academy of Sciences, Urumqi, China, in 2012. He was elected the corresponding academician from corresponding member, Royal Academy for Overseas Sciences in 2017.

He is currently a Professor with the Xinjiang Institute of Ecology and Geography, Chinese Academy of Sciences, Urumqi, China. His research interests include application of remote sensing and geographic information technology in environmental

evolution of arid areas.



Prof. Dr. Ümüt Halik studied physical geography (1982-1987, B.Sc.) and geo-ecology (1987-1990, M.Sc.) at the Xinjiang University in Urumqi, P.R. China. From 1990 to 1997, he worked as junior scientist and lecturer at the Department of Geography of the Xinjiang University. In 2002 Prof. Halik finished his PhD at the Institute for Management in Environmental Planning (IMUP) of the Technical University of Berlin (TU Berlin). His Post-Doctoral research (2003-2006) on monitoring of environmental changes and floodplain ecosystem dynamics at

the Tarim River in NW China was completed at the TU Berlin and the Xinjiang University.

Since 2006, he is full professor for landscape and urban ecology at the Xinjiang University. Within the DFG Mercator-Guest Professor-Program, Prof. Halik worked as a guest professor at the Faculty of Mathematics and Geography of Catholic University of Eichstätt-Ingolstadt in Germany (07.2010-06.2017), and held the professorship of ecosystem research for promoting international floodplain research and related networking. As project leader (PI), Prof. Halik completed several research projects founded by National Natural Science Foundation of China (NSFC), Volkswagen Stiftung, DAAD, DFG and BMBF. He has published over 100 research papers in academic magazines, and edited/published 6 monographs.

Research field: Ecosystem Ecology (Ecological Urban Planning, Watershed Ecology, Restoration Ecology), Landscape Design and Environmental Management.



Matti Pietikäinen (F12) received the D.Sc. degree in technology from the University of Oulu, Finland. He is currently a Professor, the Scientific Director of Infotech Oulu, and the Director of the Center for Machine Vision Research with the University of Oulu. From 1980 to 1981 and from 1984 to 1985, he visited the Computer Vision Laboratory, University of Maryland. He has made pioneering contributions, e.g., to local binary pattern methodology, texture-based image and video analysis, and facial image analysis. He has authored over 285 refereed papers

in international journals, books, and conferences. He was an Associate Editor of the IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE and Pattern Recognition journals, and serves as an Associate Editor of the Image and Vision Computing journal. He was the President of the Pattern Recognition Society of Finland from 1989 to 1992. From 1989 to 2007, he served as a member of the Governing Board of the International Association for Pattern Recognition (IAPR), and became one of the founding fellows of IAPR in 1994.



Wenjuan Wang received the masters degree in ecology from the school of resources and environmental sciences, Xinjiang University, Xinjiang, China, in 2019. She is currently pursuing the Ph.D. degree with the Department of Applied Ecology, Saint Petersburg State University, Saint Petersburg, Russian Federation. Her research interests include soil science and environmental toxicology.