

CNN-SVR 用于 MAX-DOAS 预测对流层 NO₂ 廓线潘屹峰¹, 田鑫^{1,2}, 谢品华^{2,3,4*}, 李昂², 徐晋², 任博^{2,4}, 黄晓辉¹, 田伟¹, 王子杰¹¹安徽大学物质科学与信息技术研究院安徽省信息材料与智能传感实验室, 安徽 合肥 230601;²中国科学院安徽光学精密机械研究所环境光学与技术重点实验室, 安徽 合肥 230031;³中国科学院区域大气环境研究卓越创新中心, 福建 厦门 361021;⁴中国科学技术大学环境科学与光电技术学院, 安徽 合肥 230026

摘要 提出一种基于卷积神经网络(CNN)和支持向量回归机(SVR)的多轴差分光学吸收光谱(MAX-DOAS)对流层 NO₂垂直分布预测方法。将 2019 年南京站点采集的原始 MAX-DOAS 数据通过 QDOAS 软件拟合获取 O₄ 和 NO₂ 差分斜柱浓度, 结合基于最优估算的气溶胶和痕量气体廓线反演算法——PriAM 算法反演了对流层 NO₂ 廓线, 并将其作为预测模型的输出。此外, 通过平均影响值方法进行预测模型输入变量的选择, 确定了 MAX-DOAS 数据、温度、气溶胶光学厚度和低云覆盖率为模型的最佳输入变量。通过实验优化网络结构和参数, 最终建立预测模型在测试集与 PriAM 的平均百分比误差仅为 9.14%, 与单独建立的 CNN、SVR、反向传播模型相比, 平均百分比误差分别降低了 8.22%、6.00%、32.28%。因此, CNN-SVR 能够利用 MAX-DOAS 数据对对流层 NO₂ 廓线进行有效预测。

关键词 大气光学; 卷积神经网络; 支持向量回归机; 多轴差分吸收光谱; 对流层 NO₂ 廓线

中图分类号 TP301

文献标志码 A

DOI: 10.3788/AOS202242.2401001

Prediction of Tropospheric NO₂ Profile Using CNN-SVR-Based MAX-DOASPan Yifeng¹, Tian Xin^{1,2}, Xie Pinhua^{2,3,4*}, Li Ang², Xu Jin², Ren Bo^{2,4}, Huang Xiaohui¹,
Tian Wei¹, Wang Zijie¹¹Information Materials and Intelligent Sensing Laboratory of Anhui Province, Institutes of Physical Science and Information Technology, Anhui University, Hefei 230601, Anhui, China;²Key Laboratory of Environmental Optics and Technology, Anhui Institute of Optics and Fine Mechanics, Chinese Academy of Sciences, Hefei 230031, Anhui, China;³CAS Center for Excellence in Regional Atmospheric Environment, Xiamen 361021, Fujian, China;⁴School of Environmental Science and Optoelectronic Technology, University of Science and Technology of China, Hefei 230026, Anhui, China

Abstract This study proposes a method based on a convolutional neural network (CNN) and support vector regression machine (SVR) for predicting the vertical distribution of NO₂ in the troposphere by multi-axis differential optical absorption spectroscopy (MAX-DOAS) technology. Taking the Nanjing site as an example, we obtain the O₄ and NO₂ differential slant column density (dSCD) according to the raw MAX-DOAS data collected by QDOAS fitting in 2019, invert the tropospheric NO₂ profile by combining the optimal estimation-based aerosol and trace gas profile inversion algorithm—PriAM, and use the profile as the output of the prediction model. In addition, the input variables of the prediction model are selected by the mean impact value method, with MAX-DOAS data, temperature, aerosol optical thickness, and low cloud coverage finally identified as the optimal input variables for the model. Furthermore, the network structure and parameters are optimized through experiments, and the average percentage error of the final CNN-SVR prediction model in the test set with PriAM is only 9.14%, which is 8.22%, 6.00%, and 32.28% lower than that of the separately constructed CNN, SVR, and backpropagation models, respectively. Therefore, CNN-SVR can effectively predict

收稿日期: 2022-03-10; 修回日期: 2022-04-17; 录用日期: 2022-05-05

基金项目: 国家自然科学基金(42105132, U19A2044, 42030609)、安徽省自然科学基金(2008085QD183, 2008085QD182)

通信作者: *phxie@aiofm.ac.cn

tropospheric NO₂ profiles by using MAX-DOAS data.

Key words atmospheric optics; convolutional neural network; support vector regression machine; multi-axis differential absorption spectroscopy; tropospheric NO₂ profile

1 引言

进入 21 世纪以来,随着人类社会经济的迅速发展,工厂以及交通排放物所占比例加大,导致以 NO₂ 为主的氮氧化物(NO_x)在大气中的浓度显著提升,使得我国面临的光化学污染和灰霾污染问题日益严峻。因此,及时掌握空气质量状况,采取更为有效的科学手段控制空气污染,对高效获取大气中 NO₂ 的空间信息提出了更为严格的要求^[1]。

多轴差分光学吸收光谱(MAX-DOAS)技术是一种广泛使用的遥测技术,可以同时观测多种大气污染成分,是一种便捷、有效地获取大气污染物垂直分布的方法^[2],已在对流层 NO₂、SO₂、HCHO 等痕量气体和气溶胶廓线反演方面得到了广泛的应用^[3-6]。现阶段利用 MAX-DOAS 技术获取痕量气体垂直分布主要基于两类反演算法,即查表法和最优估算法^[7-9],这些方法都存在一定的局限性:最优估算法依赖于先验廓线,反演过程中可能出现奇异值;查表法是一种参数化反演方法,该算法各层的分配比例和线性函数等参数定义比较粗略,空间分辨率低^[10]。同时,这两种廓线反演方法需采用大气辐射传输模型(RTM)来模拟光的传输路径,这个过程是极其复杂的,而机器学习模型可以很好地弥补廓线反演方法的这个缺点。机器学习方法通过对原始数据的分析处理和各种预测算法的应用来建立输入数据和输出数据的映射关系,直接从 MAX-DOAS 预测对流层 NO₂ 的垂直廓线,从而减少辐射传输模型模拟的流程,实现快速且准确的垂直廓线预测。

卷积神经网络(CNN)和支持向量回归机(SVR)是新颖的机器学习模型,已在光学领域得到广泛的应用。目前已有研究将 CNN 应用到近红外光谱领域进行预测和分类,均取得极佳的效果^[11-12],这也为本研究将 CNN 应用到 MAX-DOAS 数据检测中提供了理论支持。与传统的神经网络算法相比,CNN 具有更强的建模能力,可以从复杂的 MAX-DOAS 数据中提取有效的特征数据,并学习特征数据的内部结构以更好地进行预测。Malmgren-Hansen 等^[13]首次介绍了使用 CNN 从 IASI 探测数据中反演大气廓线,用深度卷积神经网络完成对多维轮廓的预测。支持向量机得益于核函数和少数起到决定作用的支持向量,适合解决各种回归预测问题。本文基于长期 MAX-DOAS 观测数据,提出一种 CNN-SVR 混合模型,并利用该模型对对流层 NO₂ 廓线进行预测,利用 CNN 和 SVR 算法建立将 MAX-DOAS 数据和气象数据等作为输入数据、对流层 NO₂ 廓线作为输出数据的模型,实现将光谱数据和气象数据等输入模型,就能得到一条准确的对流层廓线的目的。

2 实验数据与基本原理

2.1 样本数据采集

用于建模的数据来自 2019 年中国科学院安徽光学精密机械研究所(AIOFM)在南京站点(32°7'12"N, 118°57'36"E)测得的 MAX-DOAS 多仰角光谱数据。仪器方位角为 310°(规定正北方向为 0°),测量仰角时以 11 个角度为一组,分别为 1°、2°、3°、4°、5°、6°、8°、10°、15°、30°、90°。完整仰角序列的光谱保留用作本实验的输入数据,若有仰角缺失,则被判定为“有缺失”数据组,并剔除该组数据。

南京站点 MAX-DOAS 中使用的光谱仪是 Maya2000 Pro。该光谱仪的波长范围为 290~420 nm,对应 2068 条通道,光谱分辨率为 0.4~0.7 nm。将光谱仪置于恒温箱中,并将温度控制在 20 °C^[14]。本研究使用最小二乘法处理所观测的光谱^[4],利用 QDOAS(V2.111)软件进行光谱反演,NO₂ 反演的波长范围为 330~375 nm,对应的光谱数据通道数量为 613~1297 条。本实验的输入数据集即为每组测量仰角在该通道数范围内的数据。

气象数据(温度、湿度、风速、风向)对 NO₂ 的浓度有着重要影响^[15],气溶胶光学厚度(AOD)、低云覆盖率也会影响基于最优估算的气溶胶和痕量气体廓线反演算法——PriAM 算法反演对流层 NO₂ 廓线的结果^[16-18]。因此,在建立预测模型时,将这些参数也作为输入数据进行敏感性分析,求解最优敏感因子。AOD、低云覆盖率的数据来自欧洲中期天气预报中心哥白尼大气监测服务(CAMS)网站(<https://ads.atmosphere.copernicus.eu>)。目前已有大量研究证明 CAMS 预测的 AOD 的准确性,且 CAMS 相对于美国宇航局(NASA)开发的 MERRA-2 预测的中国华东地区的 AOD 准确性有相当大的提升^[19]。风向、风速、温度、相对湿度的数据来自全球天气精准预报网(<https://www.wunderground.com/history/daily/cn/nanjing/ZSNJ>)。

输出数据集是采用 MAX-DOAS 廓线反演算法 PriAM 反演的 NO₂ 廓线,具体获取流程如下:

1) NO₂ 和 O₄ 差分斜柱浓度(dSCD)的获取。基于 Lambert-Beer 定律,利用 DOAS 算法反演 NO₂ 和 O₄ 的 dSCD。将一个测量循环开始的天顶光谱作为夫琅禾费参考光谱,以去除太阳夫琅禾费结构及平流层的影响。利用 QDOAS 软件反演 NO₂ 廓线的配置和本实验选取的 330~375 nm 波段范围均是采用先前文献^[14,16-17]报道的标准配置。

2) NO₂ 廓线的反演。PriAM 算法是由 AIOFM 和 MPIC(马克斯·普朗克化学研究所)共同研发的基于非线性最优估计的痕量气体和气溶胶垂直廓线两步反

演算法。首先,将同一波段由 MAX-DOAS 反演的 O_4 dSCD 输入 PriAM,以反演气溶胶廓线。其次,利用 PriAM 痕量气体廓线反演算法,将气溶胶廓线和痕量气体 dSCD 输入到该算法中,得到对流层痕量气体垂直廓线^[14, 18-20]。

将每组原始光谱数据经过 QDOAS 拟合和 PriAM 算法反演之后得到对流层 NO_2 廓线数据。考虑到天气因素对 PriAM 算法反演对流层 NO_2 廓线的影响,剔除雨雪天气反演出的对流层 NO_2 廓线数据。剔除这部分数据的原因是雨雪天气下 PriAM 算法反演得到的对流层 NO_2 廓线数据存在很大的误差,将这些错误的输入到模型中,会影响机器学习预测模型的精度。对流层 NO_2 廓线的高度范围为 0~4 km,垂直分辨率为 200 m (最低一层为 50 m),故 NO_2 气体的浓度数据一共有 21 个格点。根据反演得到的廓线数据和对应的通道内原始光谱数据,利用神经网络算法建立预测模型。如 Hornik 等^[21]所述,只要训练数据集足够大,神经网络就可以逼近未知的函数关系,在统计学上表示为反演得到的量。本研究选取 2019 年南京站点采集的 8225 组数据作为样本,而选取连续一年的观测数据可以有效地考虑 NO_2 的季节性变化对模型验证以及误差分析的影响。

2.2 数据预处理

MAX-DOAS 数据是连续型数据,而风向、季节是离散型数据,这些参数在作为模型输入数据前,需要对离散型数据进行处理。将风向分为北(N)、东北(EN)、西北(WN)、东(E)、东南(ES)、南(S)、西南(WS)、西(W),季节分为春、夏、秋、冬。使用独热(One-Hot)编码处理离散型变量,以合理地计算特征距离^[22]。

为了提升模型精度以及收敛速度,使得本研究设计的模型能快速收敛,避免神经元饱和,对输入和输出 NO_2 廓线采用数据归一化进行预处理。归一化后的数据缩放到 [0, 1] 之间, x_i 为归一化后的数据值,归一化公式为

$$x_i = (x - x_{\min}) / (x_{\max} - x_{\min}), \quad (1)$$

式中: x 为输入数据; x_{\min} 为输入数据的最小值; x_{\max} 为输入数据的最大值。

2.3 输入变量的敏感性分析

将利用 CNN 提取的 MAX-DOAS 数据的特征向量、温度、风向、风速、相对湿度、季节、AOD、低云覆盖率作为 SVR 的输入参数。模型输入参数的好坏决定了模型的精度,本研究为了提高建模效率和实验精度,需要对输入变量进行筛选。选择使用 Dombi 等^[23]提出的平均影响值(MIV)方法,考虑到该方法原理简单、计算方便,且本研究的输入参数没有清晰的变量筛选依据,故使用应用最广泛的 MIV 方法。该方法流程如图 1 所示。

按照图 1 所示的流程,计算不同输入变量的 MIV 值占全部输入变量 MIV 之和的百分比,如图 2 所示。

MIV 的绝对值大小决定了输入变量对预测结果

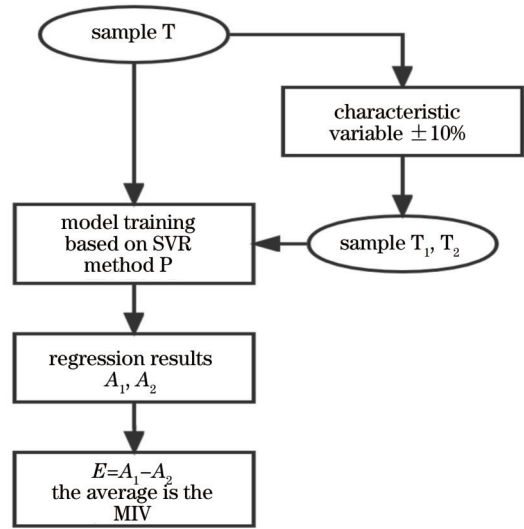


图 1 MIV 方法流程框图

Fig. 1 Flow chart of MIV method

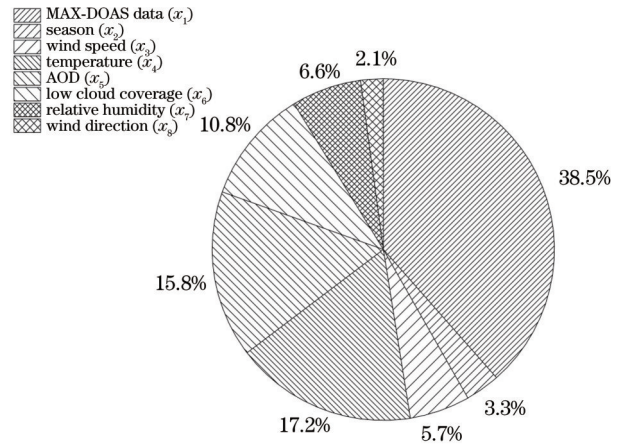


图 2 不同输入变量的 MIV 值占全部输入变量 MIV 之和的百分比

Fig. 2 Percentage of MIV value of different input variables in the sum of all input variables

的重要性,本研究得到的重要性依次为:MAX-DOAS 数据(x_1)、温度(x_4)、AOD(x_5)、低云覆盖率(x_6)、相对湿度(x_7)、风速(x_3)、季节(x_2)、风向(x_8)。根据 MIV 的绝对值从低到高的顺序依次剔除输入变量,观察不同体系中预测精度的变化,最高预测精度对应的即最优指标体系。以均方根误差(RMSE)为评价指标,其计算公式为

$$E_{\text{RMSE}} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_a - y_p)^2}, \quad (2)$$

式中: y_a 为实际值; y_p 为预测值。不同输入变量下的 RMSE 如表 1 所示。

从表 1 可以看到:输入变量体系 F1~F5 的 RMSE 值逐渐减小,说明在逐步剔除风向、季节、风速、相对湿度等变量的过程中模型的 RMSE 一直在降低;在体系 F5~F8 中逐步剔除低云覆盖率、AOD、温度,模型的 RMSE 却越来越大。当输入变量为体系 F5 时,模型

表 1 不同输入变量体系下的 RMSE 的值

Table1 RMSE values under different input variable systems

System	Input variable system	RMSE
F1	$x_1, x_4, x_5, x_6, x_7, x_3, x_2, x_8$	0.512
F2	$x_1, x_4, x_5, x_6, x_7, x_3, x_2$	0.487
F3	$x_1, x_4, x_5, x_6, x_7, x_3$	0.452
F4	x_1, x_4, x_5, x_6, x_7	0.415
F5	x_1, x_4, x_5, x_6	0.322
F6	x_1, x_4, x_5	0.401
F7	x_1, x_4	0.498
F8	x_1	0.586

RMSE 最小,即输入敏感因子中,MAX-DOAS 数据、温度、AOD 和低云覆盖率的影响明显高于风速、风向、季节和相对湿度。因此,本研究最终确定 MAX-DOAS 数据、温度、AOD 和低云覆盖率为所提模型的最佳输入变量。

2.4 CNN、SVR 和 BP 神经网络

考虑到 MAX-DOAS 的结构是一个一维向量,本研究利用 CNN 良好的自学习功能,通过 CNN 的卷积层提取 MAX-DOAS 独特的特征,池化层进行降维的同时保留 MAX-DOAS 的特征,以防止过拟合,并利用全连接层来处理得到的特征,使其输到输出层。通过 CNN 3 种不同网络层的依次作用,可以有效地提取 MAX-DOAS 的特征,避免复杂的光谱预处理。

SVR 是目前机器学习回归预测算法中的翘楚。用二元集合 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n), x_n \in \mathbf{R}, y_n \in \mathbf{R}$, 表示测试集。 x_n 为由 AOD、温度、季节、MAX-DOAS 数据、低云覆盖率等构成的特征矩阵。输出数据为对应的对流层 NO₂ 廓线 y_n 。构造和计算超平面的距离是 SVR 的核心,超平面函数 $f(x)$ 和距离 L_i 的计算公式分别为

$$f(x) = \mathbf{w}^T \mathbf{x}_n + b, \quad (3)$$

$$L_i - f(x) = L_i - \mathbf{w}^T \mathbf{x} - b, \quad (4)$$

式中: \mathbf{w}^T 为权重向量; b 为偏移量; \mathbf{x} 为输入数据的特征向量。

SVR 为了构建最优模型、得到最优超平面,引入了松弛因子 R_i 和 R'_i , 以减小回归误差,即

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + P \sum_{i=1}^n (R_i + R'_i), \quad (5)$$

$$\text{s. t. } \begin{cases} L_i - f(x) \leq \epsilon + R_i \\ f(x) - L_i \leq \epsilon + R'_i \\ R_i \geq 0, R'_i \geq 0, i = 1, 2, \dots, n \end{cases}, \quad (6)$$

式中: P 为惩罚系数,且 $P > 0$; ϵ 为不敏感损失函数。

所提 SVR 模型引入径向基函数(RBF)的核函数来解决样本数据线性不可分的问题,通过设置惩罚系数 P 和核函数系数 γ , 最终得到最优超平面的决策函数,即

$$f(x) = \sum_{i=1}^{N_s} L_i g_i K(x, x_i) + b, \quad (7)$$

式中: g_i 为支持向量; $K(\cdot)$ 为核函数; N_s 为支持向量的数量。

反向传播(BP)神经网络是应用最广泛的神经网络模型,该模型具有较强的自学习能力、泛化能力和很强的容错能力。该模型通过误差的 BP 来逐步降低模型训练中的误差,最终获得最优的 BP 神经网络模型。目前 BP 神经网络在预测近地面 NO₂ 的浓度中,基于其强大的学习和泛化能力,模型设计简单(通常使用三层 BP 神经网络),可快速预测近地面 NO₂ 的浓度,并取得不错的效果,所建立的模型具有一定的推广、概括能力。尽管 BP 神经网络目前还未曾应用于对流层 NO₂ 廓线的反演,但考虑到 BP 神经网络在预测近地面 NO₂ 的浓度取得成功,故本研究建立了 BP 神经网络模型来对比其他预测模型。

3 CNN-SVR 混合模型分析与讨论

3.1 CNN-SVR 混合模型

本研究将 CNN 和 SVR 相结合,提出一种 CNN-SVR 混合预测模型。该模型利用了 CNN 提取数据特征的能力和 SVR 在解决高维特征回归问题的高效性,将 CNN 中间层的输出特征作为提取到的最终特征并和经过 MIV 筛选的温度、AOD、低云覆盖率一起传入 SVR 模型中预测对流层 NO₂ 廓线。

首先,基于 MIV 算法的筛选确立了 MAX-DOAS 数据、温度、AOD、低云覆盖率作为模型的输入变量,输出指标数据为对流层 NO₂ 廓线。将 MAX-DOAS 数据作为输入变量输入 CNN 模型中,对原始 MAX-DOAS 数据进行归一化处理,消除量纲的影响。将 MAX-DOAS 数据输入到 CNN 卷积层,通过多次卷积提取有效特征,通过全连接层输出最终特征。以均方误差(MSE)的最小值为多次迭代的目标函数,利用梯度下降算法进行优化,从而寻找最优解,此时将全连接层的输出数据特征以矢量形式存储。将 CNN 提取的 MAX-DOAS 数据特征作为 SVR 的输入,同时将温度、AOD、低云覆盖率数据也作为 SVR 的输入,从而确定 SVR 的参数。CNN 网络层的模型参数设置如表 2 所示,所设计的 CNN-SVR 混合模型结构如图 3 所示。

为了减少由光谱平移和旋转造成的干扰,池化层采用最大池化方法采样。均方误差用来衡量 PriAM 反演的对流层 NO₂ 廓线与预测对流层 NO₂ 廓线的欧氏距离,并作为模型的损失函数。优化器使用动量梯度下降(SGD+Momentum),学习率为 10^{-3} ,批量大小(Batchsize)根据训练集的样本数量设置为 64,以保证内存利用率,加快处理速度^[24]。SVR 模型采用 RBF 核函数,为了避免陷入局部最优解,最终确定 $P=100, \gamma=0.01$ 分别作为 SVR 模型的惩罚因子和核参数^[25]。

模型及参数确定后在测试集开展对流层 NO₂ 廓线预测。采用平均绝对值误差(MAPE)[式(8)]和对称平均绝对百分比误差(SMAPE)[式(9)]来评估对

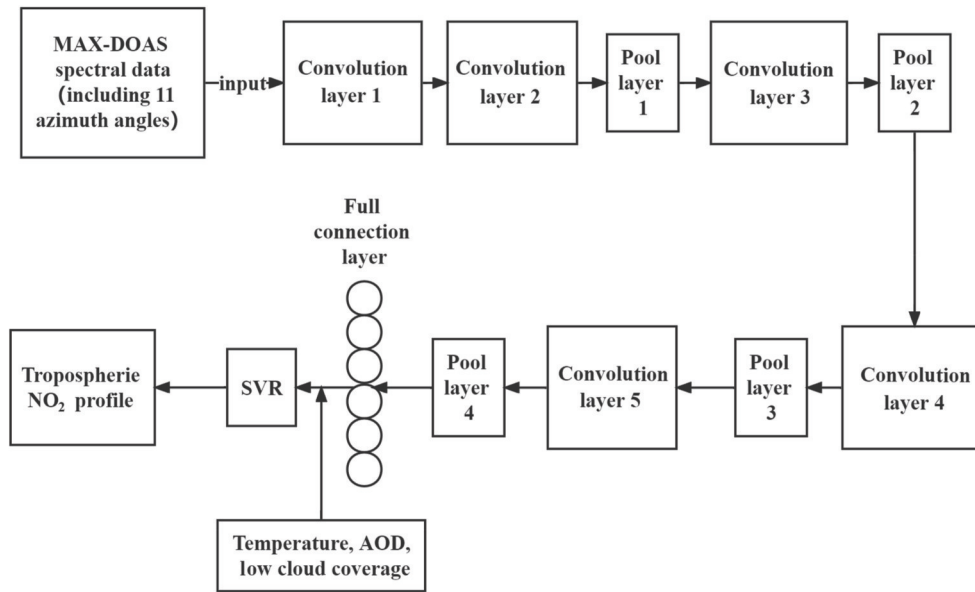


图 3 CNN-SVR混合模型框架
Fig. 3 Frame of CNN-SVR hybrid model

表 2 CNN 参数设置
Table 2 Parameter setting of CNN

Network layer	Model parameter setting
Input layer	685×11 spectral data matrix
Convolution layer 1	64 1×1 convolution kernels; kernel_size: 5
Convolution layer 2	128 1×1 convolution kernels; kernel_size: 5
Pool layer 1	MaxPool; kernel_size: 1; stride: 2
Convolution layer 3	128 1×1 convolution kernels; kernel_size: 5
Pool layer 2	MaxPool; kernel_size: 1; stride: 2
Convolution layer 4	256 1×1 convolution kernels; kernel_size: 5
Pool layer 3	MaxPool; kernel_size: 1; stride: 2
Convolution layer 5	512 1×1 convolution kernels; kernel_size: 5
Pool layer 4 (adaptive pooling layer)	Output one-dimensional vector
Full connection layer	Output 21 neurons

流层 NO₂ 廓线预测精度。MAPE 最大的优点是以百分比表示,可用于比较不同比例预测模型的准确性,易于解释预测模型的误差,但是 MAPE 的值可能超过 100% 且无上限。而 SMAPE 具有下限(0%)和上限(200%),使得误差在一定范围内而不存在无穷大值。

$$E_{MAPE} = \frac{1}{n} \sum_{i=1}^N \left| \frac{y_a - y_p}{y_p} \right| \times 100\%, \quad (8)$$

$$E_{SMAPE} = \frac{1}{n} \sum_{i=1}^N \frac{|y_a - y_p|}{(|y_a| + |y_p|)/2} \times 100\%。 \quad (9)$$

所有的原始数据样本被随机划分到训练集和测试集。为了研究 CNN-SVR 模型的不同训练集占全样本的比例对 CNN-SVR 建模效果的影响,采用相同的规则划分训练集和测试集。训练集样本占所有样本的 50%、60%、70%、80%、90%。MAPE 和 SMAPE 的测试集经过交叉验证的结果被用来评估预测的准确性。不同训练集样本下 CNN-SVR 模型测试集的预测结果如表 3 所示。当训练集的比例为 90% 时,CNN-SVR 在测试集上最佳,其 MAPE 为 9.14%,SMAPE 为 8.52%。

表 3 不同训练集样本下 CNN-SVR 模型测试集的预测结果
Table 3 Prediction results of CNN-SVR model test set under different training set samples

Ratio of training set	SMAPE / %	MAPE / %
0.5	21.97	23.41
0.6	18.11	19.33
0.7	15.13	16.24
0.8	11.55	12.41
0.9	8.52	9.14

3.2 不同预测模型对比

为了确立最优的对流层 NO₂ 廓线预测模型,对比了构建的 CNN 模型、SVR 模型、BP 神经网络和 CNN-SVR 混合模型。

从表 4 所示的 4 种模型(CNN-SVR、BP、CNN、SVR)在训练数据集和测试数据集上的 MAPE 和

SMAPE 可以看出, CNN-SVR 混合模型的预测结果误差最小。与 CNN 模型、SVR 模型、BP 模型相比, CNN-SVR 混合模型的 MAPE 分别降低了 8.22%、6.00%、32.28%, SMAPE 分别降低了 7.59%、5.76%、30.69%。因此, 使用 CNN-SVR 混合模型预测对流层 NO₂ 廓线的效果最优, 具有更高的预测精度。

表 4 4 种预测模型下的 MAPE 和 SMAPE
Table 4 MAPE and SMAPE under four prediction models

Model	MAPE / %		SMAPE / %	
	Training error	Test error	Training error	Test error
CNN	12.59	17.63	11.61	16.11
SVR	11.82	15.14	10.92	14.28
BP	35.68	41.42	33.63	39.21
CNN-SVR	7.93	9.14	7.25	8.52

4 种模型 (CNN-SVR、BP、CNN、SVR) 的训练集和测试集预测的 NO₂ 廓线与 PriAM 反演的真实 NO₂ 廓线之间的决定系数 R² 如表 5 所示。CNN-SVR 模型的决定系数较其他 3 种模型的结果高, CNN-SVR 模型

在训练集的 R² 为 0.95, 在测试集的 R² 为 0.93, 表明 CNN-SVR 拟合优度最佳。

表 5 4 种预测模型下的 R²
Table 5 R² under four prediction models

Model	R ² in modeling set	R ² in testing set
CNN	0.86	0.83
SVR	0.89	0.87
BP	0.81	0.77
CNN-SVR	0.95	0.93

3.3 对比验证

为了验证建立的 CNN-SVR 模型的预测能力, 在春、夏、秋、冬 4 个季节中随机各取 1 个晴天, 利用 CNN-SVR 模型预测这 4 天的 NO₂ 廓线, 并与 PriAM 反演的 NO₂ 廓线进行对比, 结果如图 4 所示。总体来说, CNN-SVR 预测的对流层 NO₂ 廓线与 PriAM 反演的 NO₂ 廓线基本一致, 且时间序列也高度一致, 说明 CNN-SVR 模型可以很好地预测对流层 NO₂ 廓线, 其预测的近地面 NO₂ 浓度值与 PriAM 反演结果最大相差约为 15%, 在 0.6 km 以上相较于 PriAM 几乎没有

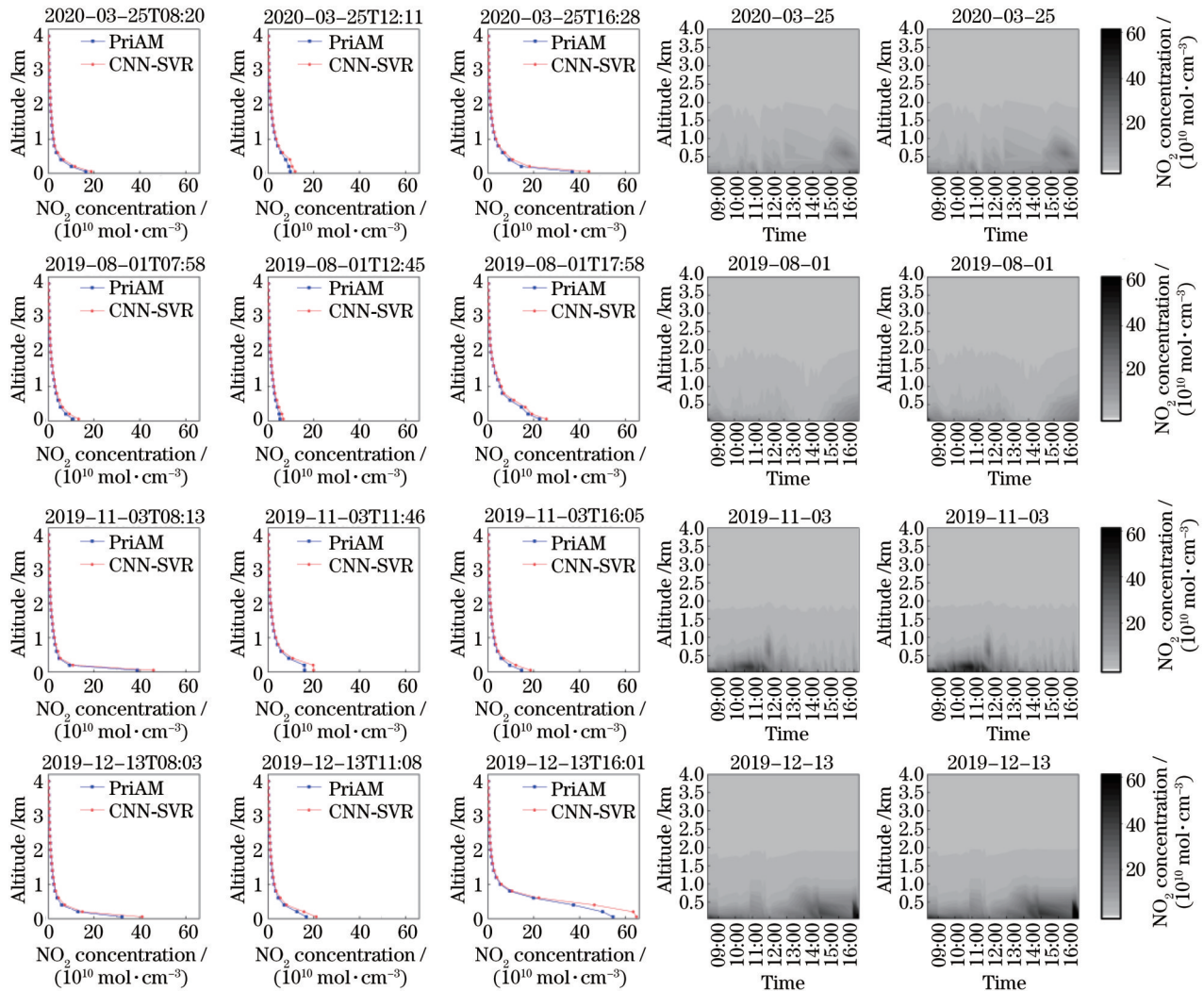


图 4 CNN-SVR 模型预测与 MAX-DOAS 反演的 NO₂ 廓线对比

Fig. 4 Comparison of NO₂ profiles predicted by CNN-SVR model and retrieved by MAX-DOAS

差异。

为了进一步对比算法的可靠性和优劣性,利用 2019 年 7 月至 2020 年 6 月的 MAX-DOAS 数据对其 NO_2 廓线进行预测。选择该时段数据的主要原因是: 1) 2019 年 7—12 月数据包括预测模型的建模集数据, 可能导致预测数据的误差有一定偏向性, 因此再利用 2020 年 1—6 月的月平均廓线计算其误差来验证模型的普适性; 2) 对流层 NO_2 的浓度受到季节因素的影响, 使用连续一年的数据可以更好地验证模型, 避免了一季节的因素导致预测数据误差偏高或偏低。CNN-

SVR 模型预测的以及 PriAM 反演的 NO_2 的月平均廓线对比如图 5 所示。利用两种方法得到的 NO_2 廓线具有较好的一致性, 表明 CNN-SVR 模型可以很好地预测对流层 NO_2 廓线。总体上 CNN-SVR 模型预测的 NO_2 浓度较 PriAM 反演的结果高约 9.4%。反演误差主要集中在 1 km 及以下, 该区域廓线的平均 MAPE 为 18%。这是因为输出的是 21 个不同高度的对流层 NO_2 廓线, CNN-SVR 预测模型在回归的过程中为了维持总体的误差, 保证绝大多数高度可以预测准确, 在预测中没有偏向将 1 km 以下的数据预测准确。

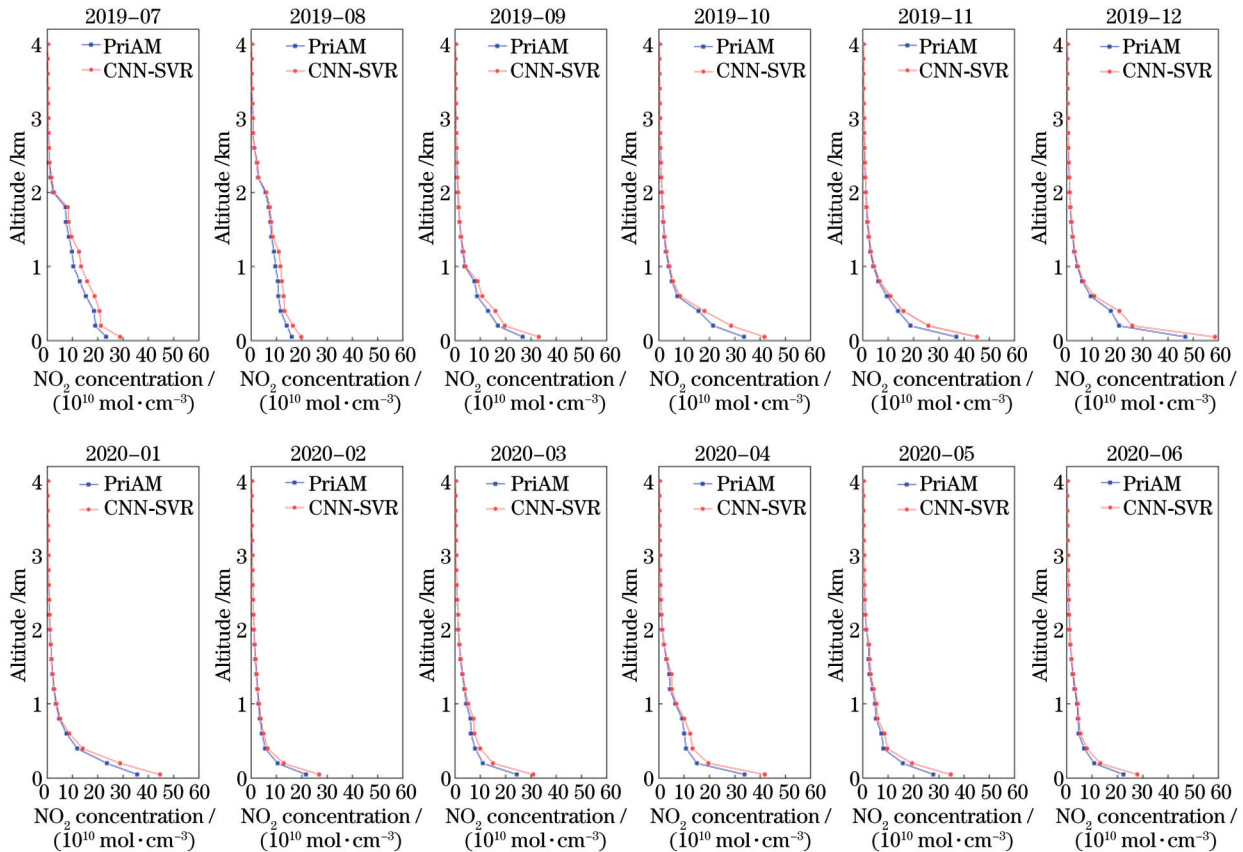


图 5 2019 年 7 月至 2020 年 6 月的 NO_2 月平均廓线

Fig. 5 Monthly average profile of NO_2 from July 2019 to June 2020

3.4 不同 AOD 的敏感性分析

由于气溶胶对于痕量气体廓线反演的影响较大, AOD 的增大会一定程度地增加 PriAM 算法的平滑误差。研究发现, 气溶胶 AOD 的错误对 NO_2 反演的影响尤为剧烈, 气溶胶 AOD 的增加会在一定程度降低反演廓线的灵敏度, 因此本研究开展了 AOD 对于 CNN-SVR 模型预测敏感性的分析。从 2019 年 7 月至 2020 年 6 月的数据也发现了在不同的 AOD 条件下, CNN-SVR 预测结果与 PriAM 反演数据的平均百分比误差存在一定差异, 结果如图 6 所示。当 AOD 大于 3.0 时, CNN-SVR 与 PriAM 的平均 MAPE 最大 (12.08%); 当 AOD 在 0~0.5 时, 平均 MAPE 最小 (6.72%)。不同 AOD 下 CNN-SVR 模型预测的 MAPE 范围如表 6 所示。可以明显看出, AOD 在 0~0.5 范围内, 预测

MAPE 的最大值为 7.63%, 而 AOD 大于 3.0 的 MAPE 最小值高达 9.71%, 最大预测 MAPE 更是 AOD 范围在 0~0.5 的 2 倍多。表 6 和图 6 的数据进一步证明了 AOD 的大小会影响 PriAM 反演的准确性, 这是因为 CNN-SVR 模型的训练集中包含了一定数量的、由 AOD 的影响造成 PriAM 反演存在误差的数据, 这些数据会对 CNN-SVR 模型最终的结果造成负面影响, 从而降低模型精度。

4 结 论

将 CNN 和 SVR 用于 MAX-DOAS 预测对流层 NO_2 廓线, 提出一种有效的 CNN-SVR 模型, 实现了通过 MAX-DOAS 原始光谱数据来实值回归预测对流层 NO_2 廓线, 并取得了较好效果。通过 MAX-DOAS 数

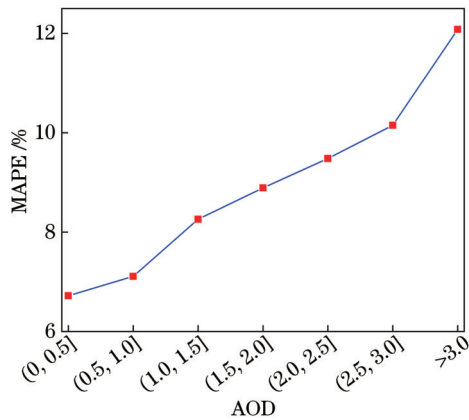


图6 不同AOD范围内CNN-SVR模型预测误差

Fig. 6 Prediction error of CNN-SVR model under different AOD ranges

表6 不同AOD范围内CNN-SVR模型预测误差的范围

Table 6 Prediction error range of CNN-SVR model under different AOD ranges

AOD	Prediction MAPE range / %
(0, 0.5]	(3.17, 7.63)
(0.5, 1.0]	(5.13, 8.21)
(1.0, 1.5]	(6.73, 9.16)
(1.5, 2.0]	(6.15, 11.19)
(2.0, 2.5]	(7.59, 11.80)
(2.5, 3.0]	(7.21, 13.17)
>3.0	(9.71, 16.23)

据、温度、AOD、低云覆盖率和PriAM反演的真实对流层NO₂廓线进行模拟实验,实验结果证明了该模型的优越性。CNN可以提取更复杂的特征,提高计算速度;采用SVR算法对提取的特征进行预测,可提高预测精度。结果表明,与CNN模型、SVR模型和BP模型相比,CNN-SVR混合模型的误差减小了8%~32%,能够快速准确地预测目标值,性能更加稳定。同时,选取了小时廓线、日均廓线和月均廓线结果进行对比,发现CNN模型预测的廓线更接近MAX-DOAS反演的真实NO₂廓线,说明了将MAX-DOAS数据输入CNN-SVR模型预测NO₂廓线的方法是有效的,尤其是CNN可以从原始光谱中学习关键模式,显著降低对特征工程的需求。多个处理层的组合提高了拟合和特征提取能力,使其适用于学习MAX-DOAS数据潜在特征,利用MAX-DOAS数据特征进行对流层NO₂廓线的实值回归预测。不同季节随机几天的对比结果表明CNN-SVR模型能较好地预测对流层NO₂廓线,预测结果与PriAM反演结果高度一致。利用建立的CNN-SVR模型对2019年7月至2020年6月南京站点的NO₂廓线进行预测,发现CNN-SVR模型能较好地预测对流层NO₂廓线,是一种有效的对流层NO₂廓线预测方法。同时分析了不同AOD大小对CNN-SVR模型预测精度的影响,AOD越大,模型预测精度越低。

如今机器学习模型对于各种应用均取得不错的效

果,但是在模型的设置中需要相当多的技能和经验来选择合适的参数,极度依赖精确的训练集输出数据和合适的输入数据。为了更准确地预测对流层NO₂廓线,需要找到更为适合的输入因子,这是预测对流层NO₂廓线的难点,也是未来基于MAX-DOAS数据通过CNN-SVR模型预测对流层NO₂廓线可进一步开展的工作。

参 考 文 献

- [1] 周海金,刘文清,司福祺,等. 被动多轴差分吸收光谱技术监测大气NO₂垂直廓线研究[J]. 光学学报, 2011, 31(11): 1101007.
Zhou H J, Liu W Q, Si F Q, et al. Retrieval of atmospheric NO₂ vertical profile from multi-axis differential optical absorption spectroscopy[J]. Acta Optica Sinica, 2011, 31(11): 1101007.
- [2] Platt U, Stutz J. Differential absorption spectroscopy [M]//Platt U, Stutz J. Differential optical absorption spectroscopy. Physics of earth and space environments. Heidelberg: Springer, 2008: 135-174.
- [3] Frieß U, Sihler H, Sander R, et al. The vertical distribution of BrO and aerosols in the Arctic: measurements by active and passive differential optical absorption spectroscopy[J]. Journal of Geophysical Research: Atmospheres, 2011, 116(D14): D00R04.
- [4] Hendrick F, Müller J F, Clémer K, et al. Four years of ground-based MAX-DOAS observations of HONO and NO₂ in the Beijing area[J]. Atmospheric Chemistry and Physics, 2014, 14(2): 765-781.
- [5] Kanaya Y, Irie H, Takashima H, et al. Long-term MAX-DOAS network observations of NO₂ in Russia and Asia (MADRAS) during the period 2007-2012: instrumentation, elucidation of climatology, and comparisons with OMI satellite observations and global model simulations[J]. Atmospheric Chemistry and Physics, 2014, 14(15): 7909-7927.
- [6] Li X, Brauers T, Hofzumahaus A, et al. MAX-DOAS measurements of NO₂, HCHO and CHOCHO at a rural site in Southern China[J]. Atmospheric Chemistry and Physics, 2013, 13(4): 2133-2151.
- [7] Hönninger G, von Friedeburg C, Platt U. Multi axis differential optical absorption spectroscopy (MAX-DOAS)[J]. Atmospheric Chemistry and Physics, 2004, 4 (1): 231-254.
- [8] 吴丰成,谢品华,李昂,等. 基于多轴差分吸收光谱技术的查找表法反演气溶胶消光廓线研究[J]. 光学学报, 2013, 33(6): 0601002.
Wu F C, Xie P H, Li A, et al. Research of aerosol extinction inverted with look-up table method based on multi-axis differential optical absorption spectroscopy[J]. Acta Optica Sinica, 2013, 33(6): 0601002.
- [9] 刘文清. “双碳”目标下大气环境光学监测技术发展机遇[J]. 光学学报, 2022, 42(6): 0600001.
Liu W Q. Opportunities and challenges for development of atmospheric environmental optics monitoring technique under "double carbon" goal[J]. Acta Optica Sinica, 2022,

- 42(6): 0600001.
- [10] Irie H, Takashima H, Kanaya Y, et al. Eight-component retrievals from ground-based MAX-DOAS observations[J]. *Atmospheric Measurement Techniques*, 2011, 4(6): 1027-1044.
- [11] 王璨, 武新慧, 李恋卿, 等. 卷积神经网络用于近红外光谱预测土壤含水率[J]. *光谱学与光谱分析*, 2018, 38(1): 36-41.
- Wang C, Wu X H, Li L Q, et al. Convolutional neural network application in prediction of soil moisture content [J]. *Spectroscopy and Spectral Analysis*, 2018, 38(1): 36-41.
- [12] 唐永生, 陈争光. 卷积神经网络和近红外光谱的土壤 pH 值预测 [J]. *光谱学与光谱分析*, 2021, 41(3): 892-897.
- Tang Y S, Chen Z G. Soil pH prediction based on convolution neural network and near infrared spectroscopy [J]. *Spectroscopy and Spectral Analysis*, 2021, 41(3): 892-897.
- [13] Malmgren-Hansen D, Nielsen A A, Laparra V, et al. Transfer learning with convolutional networks for atmospheric parameter retrieval[C]//IGARSS 2018-2018 IEEE International Geoscience and Remote Sensing Symposium, July 22-27, 2018, Valencia, Spain. New York: IEEE Press, 2018: 2111-2114.
- [14] Tian X, Xie P H, Xu J, et al. Long-term observations of tropospheric NO₂, SO₂ and HCHO by MAX-DOAS in Yangtze River Delta area, China[J]. *Journal of Environmental Sciences*, 2018, 71: 207-221.
- [15] 熊险平, 何璇, 赵玉广. NO₂浓度变化及气象因子清除能力分析[J]. *环境科学导刊*, 2020, 39(4): 27-33.
- Xiong X P, He X, Zhao Y G. The analysis of NO₂ concentration changes and removal ability of meteorological factors[J]. *Environmental Science Survey*, 2020, 39(4): 27-33.
- [16] Ren B, Xie P H, Xu J, et al. Use of the PSCF method to analyze the variations of potential sources and transports of NO₂, SO₂, and HCHO observed by MAX-DOAS in Nanjing, China during 2019[J]. *Science of the Total Environment*, 2021, 782: 146865.
- [17] Wang Y, Lampel J, Xie P H, et al. Ground-based MAX-DOAS observations of tropospheric aerosols, NO₂, SO₂ and HCHO in Wuxi, China, from 2011 to 2014[J]. *Atmospheric Chemistry and Physics*, 2017, 17(3): 2189-2215.
- [18] 王杨, 李昂, 谢品华, 等. 多轴差分吸收光谱技术测量 NO₂ 对流层垂直分布及垂直柱浓度 [J]. *物理学报*, 2013, 62(20): 200705.
- Wang Y, Li A, Xie P H, et al. Measuring tropospheric vertical distribution and vertical column density of NO₂ by multi-axis differential optical absorption spectroscopy[J]. *Acta Physica Sinica*, 2013, 62(20): 200705.
- [19] Fu D S, Liu M Q, Yang D Z, et al. Influences of atmospheric reanalysis on the accuracy of clear-sky irradiance estimates: comparing MERRA-2 and CAMS [J]. *Atmospheric Environment*, 2022, 277: 119080.
- [20] 田鑫, 徐晋, 谢品华, 等. 基于多轴差分吸收光谱技术测量对流层 HCHO 垂直分布 [J]. *光谱学与光谱分析*, 2019, 39(8): 2325-2331.
- Tian X, Xu J, Xie P H, et al. Retrieving tropospheric vertical distribution in HCHO by multi-axis differential optical absorption spectroscopy[J]. *Spectroscopy and Spectral Analysis*, 2019, 39(8): 2325-2331.
- [21] Hornik K, Stinchcombe M, White H. Multilayer feedforward networks are universal approximators[J]. *Neural Networks*, 1989, 2(5): 359-366.
- [22] Buckman J, Roy A, Raffel C, et al. Thermometer encoding: one hot way to resist adversarial examples[C]//6th International Conference on Learning Representations 2018, April 30-May 3, 2018, Vancouver, BC, Canada. [S.l.: s.n.], 2018.
- [23] Dombi G W, Nandi P, Saxe J M, et al. Prediction of rib fracture injury outcome by an artificial neural network[J]. *The Journal of Trauma*, 1995, 39(5): 915-921.
- [24] 张祥东, 王腾军, 朱劭俊, 等. 基于扩张卷积注意力神经网络的高光谱图像分类 [J]. *光学学报*, 2021, 41(3): 0310001.
- Zhang X D, Wang T J, Zhu S J, et al. Hyperspectral image classification based on dilated convolutional attention neural network[J]. *Acta Optica Sinica*, 2021, 41(3): 0310001.
- [25] 朱黎明, 孙刚, 陈多龙, 等. 基于支持向量机估算大气光学湍流廓线的研究 [J]. *光学学报*, 2022, 42(1): 0101001.
- Zhu L M, Sun G, Chen D L, et al. Atmospheric optical turbulence profile estimation using support vector machine [J]. *Acta Optica Sinica*, 2022, 42(1): 0101001.