

文章编号: 2095-4980(2023)10-1247-10

## 基于红外漫反射谱和机器学习的粉末物质识别

高 颂<sup>1</sup>, 阎结昀<sup>1</sup>, 王迎新<sup>2</sup>, 解 研<sup>2</sup>

(1.北京邮电大学 理学院, 北京 100876; 2.清华大学 工程物理系, 北京 100084)

**摘 要:** 红外光谱可有效携带化合物结构以及化合物组成成分的信息, 在化学研究、纯度检测和药物识别领域已得到广泛应用。但在实际应用场景中, 由于缺乏标准样品制备的条件, 红外光谱识别准确率较低, 效率较差, 使这项技术受到了极大限制。本文采用可调谐红外量子级联激光器作为激光源, 建立记录粉末样品漫反射光谱的实验系统。以葡萄糖和聚乙烯混合粉末漫反射谱为例, 通过 Kubelka-Munk(K-M)方程和 Kramers-Kronig(K-K)关系合成样品的透射谱, 并验证漫反射谱还原透射谱的可能性。将光谱数据用于两种神经网络模型中, 对混合粉末质量分数进行预测。结果表明, 在 K-K 关系变换下, 长短期记忆(LSTM)网络模型预测效果最佳, 明显优于 BP 神经网络模型; 在 K-M 方程变换下, 两种神经网络对高质量分数葡萄糖样品的预测都比较准确, 对低质量分数的预测较差。两种漫反射光谱校正方法都不同程度地提高了训练结果的准确性, LSTM 模型整体优于 BP 神经网络模型。这些研究结果有助于发展基于频率可调谐或宽谱红外激光的未知混合粉末样品的识别技术。

**关键词:** 量子级联激光器; 漫反射; 红外光谱法; Kramers-Kronig 关系; 人工神经网络

中图分类号: TN21

文献标志码: A

doi: 10.11805/TKYDA2022035

## Powder compound identification based on infrared diffuse reflectance spectroscopy and machine learning

GAO Song<sup>1</sup>, YAN Jieyun<sup>1</sup>, WANG Yingxin<sup>2</sup>, XIE Yan<sup>2</sup>

(1.School of Science, Beijing University of Posts and Telecommunications, Beijing 100876, China;

2.Department of Engineering Physics, Tsinghua University, Beijing 100084, China)

**Abstract:** Infrared spectroscopy can effectively carry the information of compound structure and compound composition, which has been widely used in chemical research, purity detection and drug identification. However, in practical application scenarios, due to the lack of standard sample preparation conditions, the low accuracy and low recognition efficiency of infrared spectroscopy also make this technology limited. By using a tunable infrared Quantum Cascade Laser(QCL) as the source, an experimental system for recording the diffuse reflection spectra of powder samples is established. Taking the diffuse reflectance spectrum of a mixed powder of glucose and polyethylene as an example, the Kubelka-Munk(K-M) equation and the Kramers-Kronig(K-K) relations are utilized to synthesize the transmission spectrum of the sample from the experimentally measured diffuse reflectance spectrum of pure glucose, and the possibility of the transformation of the transmission spectrum from the diffuse reflectance spectrum is verified. On this basis, spectroscopic data are applied to two neural network models to predict the mass fraction of mixed powder. The results show that the Long Short-Term Memory (LSTM) predicts the best results under the K-K relations, significantly better than the BP neural network. Under the K-M equation, both neural networks are more accurate in predicting the glucose samples with high mass fraction and poorer in predicting the glucose samples with low mass fraction. Both diffuse reflectance spectral correction methods improve the accuracy of the training results, and the LSTM prefers to the BP neural network. This work contributes to the development of the identification of unknown mixed powder samples based on frequency-tunable or broad spectral infrared lasers.

**Keywords:** Quantum Cascade Laser; diffuse reflection; infrared spectroscopy; Kramers-Kronig relations; artificial neural network

太赫兹波是指频率在 0.1~10 THz 的电磁波, 波长在 0.03~3 mm 之间, 是一种红外波。太赫兹波在物质识别、光谱分析中有着独特的优势, 在大气成分监测、痕量物质监测等领域具有重要的应用前景<sup>[1]</sup>。中红外激光是指波长在 2.5~25  $\mu\text{m}$  之间的电磁波, 频率在 12~120 THz 之间, 近年来亦被称为深太赫兹波段而受到太赫兹科学和技术领域研究人员的关注。因其在空气中传播性较强, 在物质中透过性较好, 当辐照在物质表面时, 可引起分子的收缩振动, 使分子内电子发生能级间的跃迁现象, 可反映出物质中化学键和官能团等信息。由于光谱信息的稳定性取决于待测物质化学成分的稳定性, 当物质成分发生稳定变化时, 这些变化在红外光谱中也有相应的强度变化关系。基于这些特点, 利用中红外进行物质识别和对物质光谱深入研究一直是领域内的热门话题<sup>[2-4]</sup>。

红外光谱法(Infrared Spectroscopy, IR)是一种能够有效研究化合物结构以及组成成分的光学方法, 是分子吸收光谱的一种, 其作为化学研究中的一种重要鉴别和分析方法, 已得到广泛应用<sup>[5]</sup>。对于不同聚集状态的物质, 通常采用傅里叶变换红外光谱(Fourier Transform Infrared Spectroscopy, FTIR)的方法来分析化合物性质, 这种方法大量用于气体化合物的分析工作。目前对于固体和液体的化合物分析主要采用拉曼光谱方法<sup>[6]</sup>。但在常温状态下, 分子获得的能量较小, 大部分电子处于未激发态, 不能有效反映物质的结构信息。若对化合物进行定量分析, 需对样品的制备以及实验环境有一定的要求<sup>[7]</sup>。

太赫兹时域光谱技术是太赫兹光谱技术研究的主要方法, 和 FTIR 同样用于中红外波段以及远红外波段。太赫兹时域光谱技术通过时间延迟的方式获取时域电场强度的物理量, 然后通过傅里叶变换从时域信息中获取频域中的光谱信息<sup>[8-9]</sup>。由于太赫兹时域光谱技术通过脉冲的方式激发太赫兹脉冲信号, 不同于 FTIR 采用红外热辐射的方式产生光源, 因此这种同步相干的探测方式使太赫兹时域光谱技术受热辐射背景噪声的影响较小。但在实际应用中, 为保证高性能、高精度, 常采用制冷型太赫兹探测技术, 太赫兹探测器设计较为复杂, 体积庞大, 如果在室温环境下, 探测器的灵敏度和带宽往往达不到分析要求<sup>[10]</sup>。若用外差法进行探测, 需进行大量的准备工作建立实验系统, 且探测效率低, 目前尚不能满足投入实际应用的需求。而中红外光谱分析技术因其工作频率范围大, 解决了光谱带宽的限制, 同时室温工作的连续波激光光源及室温工作的探测器产品已经成熟, 集成度高等特点也一定程度上解决了太赫兹漫反射光谱研究工作中的难题。

1994 年, 美国贝尔实验室首次在超晶格结构中实现了能够发射中红外和远红外频段激光的量子级联激光器(QCL)。此后, 可调谐红外量子级联激光器被成功研制, 其可调谐范围达到  $1\ 000\ \text{cm}^{-1}$ , 峰值功率高达 150 mW<sup>[11]</sup>。实际实验环境中, 漫反射光谱辐射较弱, 对比吸收光谱选择性差, 给研究工作带来了巨大困难, 这种高功率宽谱可调谐红外激光器的发明给漫反射光谱分析带来了可能。蒋涛等采用固源分子束外延技术, 基于半绝缘-等离子体波导工艺制作了太赫兹量子级联激光器, 其在 10 K 工作温度、350 mA 激励电流下的中心频率为 2.93 THz, 电流密度达到  $156\ \text{A}/\text{cm}^2$ , 输出功率达到 7.84 mW, 提高了太赫兹器件的性能<sup>[12]</sup>。Erik R Deutsch 等利用红外 QCL 实现了表面痕量污染物的漫反射光谱检测<sup>[13]</sup>。Paul Dean 等<sup>[14]</sup>和 I S Golyak 等<sup>[15]</sup>分别通过太赫兹 QCL 和中红外 QCL 研究了在粉末样品漫反射光谱中获取有效吸收系数的方法, 证明了漫反射测量在粉末样品分析中的适用性。光谱分析与常规的化学方法制样分析不同, 它不能通过化学反应有损地对物质进行鉴定, 因此常采用计算机技术中的机器学习方法来建立模型。通过对模型的参数约束, 光谱中携带的吸收特征信息可从重叠的光谱数据集中被提取出来, 从而反映出待测物质的分子结构、性质等物理参数和化学参数, 也可以在检测混合物纯度、成分中加以应用。张活等<sup>[16]</sup>基于太赫兹时域光谱技术建立了最小二乘、区间最小二乘等回归模型, 预测了三七粉中掺杂的淀粉成分。张航等<sup>[17]</sup>基于傅里叶红外光谱技术建立了最小二乘回归模型, 并引入了遗传算法、竞争自适应重加权采样等特征选择方法, 实现了乙醇固态发酵过程的参数预测。L Wang<sup>[18]</sup>等基于傅里叶红外光谱技术建立了支持向量机和反向传播神经网络模型, 识别出皂荚样品中掺杂的蔷薇和玫瑰成分。黄妙芬等<sup>[19]</sup>利用卫星采集的可见光-近红外光谱数据建立了长短期记忆(LSTM)网络模型, 预测了海水中石油污染的含量。

本文采用商用可调谐红外 QCL 作为激光源, 记录葡萄糖粉末和葡萄糖、聚乙烯在不同浓度下的混合粉末漫反射光谱。通过 Kubelka-Munk(K-M)方程和 Kramers-Kronig(K-K)方程分别计算样品的有效吸收系数, 进一步验证两种漫反射光谱校正方法的适用性, 并提出了基于漫反射光谱校正的预处理方法, 通过构建 BP 神经网络、LSTM 模型, 实现了混合粉末样品的浓度识别。

## 1 实验系统和实验样品

### 1.1 基于频率可调谐红外 QCL 的粉末样品漫反射谱实验系统

本文采用的激光器为 Block Engineering(美国)公司生产的 LaserTune 可调谐红外量子级联激光器, 工作波段在  $6.09\sim 10.61\ \mu\text{m}$ , 峰值功率约为  $101.2\ \text{mW}$ , 平均功率变化范围在  $0.52\sim 2.52\ \text{mW}$ , 激光器和光电探测器以脉冲的模式实现采集信息的同步。为消除实验光路长度以及测试样品对激光器与光电探测器同步率的影响, 本实验中将激光器脉冲持续时间设置为  $60\ \text{ns}$ , 重复频率为  $0.5\ \text{MHz}$ ; 光电探测器采样持续时间设置为  $70\ \text{ns}$ , 采样延迟根据光路长度等因素设置为  $110\ \text{ns}$ 。同时为了获取较高精确度的漫反射光谱, 激光调谐步长设置为  $0.5\ \text{cm}^{-1}$ 。

实验光路如图 1 所示, 其中 1 为可调谐红外 QCL; 2 为碲镉汞(MCT)光电探测器; 3 为测试样本; 4 为平面镜; 5 为  $45^\circ$  离轴抛物面镜; 6 为锗透镜; 7 为 QCL-MCT 光电探测器连接线; 8 为 QCL-PC 连接线; 9 为 PC; 1-3 间引导线为入射辐射; 3-2 间引导线为漫反射辐射。QCL 工作产生的红外辐射通过平面镜以及  $45^\circ$  离轴抛物面镜反射聚焦在测试样品上, 产生吸收、反射和漫反射现象。将样品位置替换为平面镜, 将光电探测器置于反射角位置, 如收集反射辐射得到的功率谱与激光器出厂参考报告功率谱有强相关性, 说明测量光路具有良好的可靠性。为了收集到漫反射辐射且受镜面反射辐射影响最小, 将光电探测器置于法线位置, 并由一个 Ge 透镜收集漫反射辐射。测量数据通过两条连接线同步传输到 PC 中, 整个光路的长度大约在  $50\ \text{cm}$  左右。

### 1.2 测试样本数据获取

在实验数据获取阶段, 为消除激光器功率波动对结果的影响, 首先将探测器置于图 1 中 1 和 4 之间, 将测得激光器功率波动范围较小的  $942.75\sim 1300.75\ \text{cm}^{-1}$  作为研究波段, 并记录结果作为参考激光器功率信号  $I_0$ 。实验所用的聚乙烯样品是颗粒大小为  $56\sim 75\ \mu\text{m}$  的白色粉末, 无需特殊制备; 葡萄糖粉末颗粒较大, 需在玛瑙研钵中进行  $2\sim 3\ \text{min}$  的研磨制备, 将粉末研磨至  $60\ \mu\text{m}$  左右。然后将测试样品置于  $5\ \text{cm}\times 5\ \text{cm}\times 2\ \text{cm}$  聚丙烯方形容器内并填满, 放置于图 1 中 3 的位置。

将样品整理平整, 以减少入射辐射在样品内部发生多次反射吸收现象, 获得较高强度的漫反射辐射信号。测量葡萄糖粉末样品光谱 20 组, 以及葡萄糖和聚乙烯混合粉末在 6 个不同质量分数下(聚乙烯质量分数分别为  $0.01$ 、 $0.05$ 、 $0.1$ 、 $0.3$ 、 $0.5$ 、 $0.7$ ), 不同点位的实验光谱, 每组质量分数下测量 20 组, 共 120 组, 每组数据有 731 个数据点, 记为  $I_1$ 。计算  $R=I_1/I_0$  作为漫反射率, 混合粉末样本数据集如图 2 所示, 横坐标为波数, 纵坐标为漫反射率。

## 2 红外漫反射光谱分析理论

### 2.1 Kubelka-Munk 理论

研究者们建立了许多理论模型来描述光在介质中发生漫反射辐射的现象, 尤其是光在非均匀介质中的传播。其中最常用的是 Kubelka-Munk 理论, 因为它的公式化理论较为简单, 通过对辐射传输方程的推导, 描述了均匀漫反射辐射在一维半无限各向同性空间的传播过程<sup>[20]</sup>。

$$F(R_\infty) = \frac{k}{s} = \frac{(1-R_\infty)^2}{2R_\infty} \quad (1)$$

式中:  $R_\infty$  为样品在无限厚时的反射率;  $k$  为样品的吸收系数;  $s$  为样品的散射系数。实际应用中, 样品厚度并不需要无限厚, 一般在几个毫米以上的就可以适用 K-M 方程。

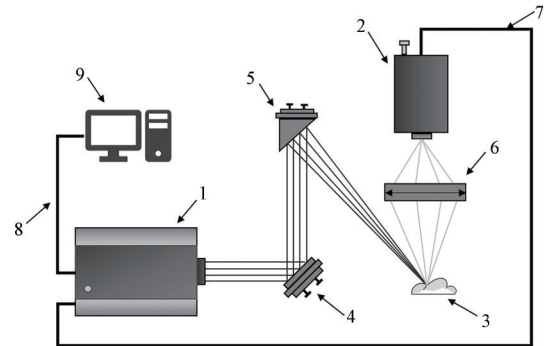


Fig.1 Device for collecting diffuse reflectance signals from solid samples

图 1 收集固体样品漫反射实验装置

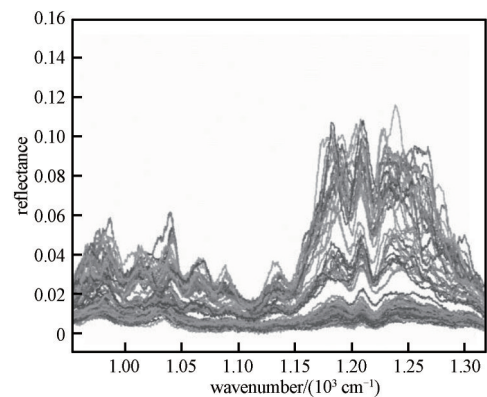


Fig.2 Diffuse reflectance spectra of mixed glucose/PE powder  
图 2 葡萄糖和葡萄糖/聚乙烯混合粉末漫反射光谱曲线图

## 2.2 Kramers-Kronig 理论

当光入射至两种介质的界面时, 根据菲涅耳方程, 反射、折射和吸收系数之间存在代数关系, 介质的反射率可以表示为一个复函数:

$$\hat{r}(\nu) = \frac{(n-1) + i\alpha}{(n+1) + i\alpha} = \eta e^{i\phi} \quad (2)$$

式中:  $n$  为样品的折射率;  $\alpha$  为样品的消光系数;  $\eta$  为反射波的振幅。本文中, 光通过空气和样品两种介质。

对等式两端取对数得  $\ln[\hat{r}(\nu)] = \ln[\eta(\nu)] + i\phi(\nu)$ , 函数  $\ln[\eta(\nu)]$  和  $\phi(\nu)$  为复函数的实部和虚部, 它们之间满足 Kramers-Kronig 关系<sup>[21]</sup>:

$$\ln[\eta(\nu)] = P \left( \frac{2}{\pi} \int_0^{\infty} \frac{\nu' \phi(\nu')}{\nu'^2 - \nu^2} d\nu' \right) \quad (3)$$

$$\phi(\nu) = P \left( - \frac{2\nu}{\pi} \int_0^{\infty} \frac{\ln[\eta(\nu')]}{\nu'^2 - \nu^2} d\nu' \right) \quad (4)$$

式中:  $P$  为主值积分;  $\nu'$  为吸收峰位置的振动频率。

由于漫反射辐射  $R(\nu) = \eta^2(\nu)$ , 可以通过式(5)来计算吸收系数  $k$ :

$$k(\nu) = \frac{2\eta(\nu)\sin\phi(\nu)}{1 - 2\eta(\nu)\cos\phi(\nu) + R(\nu)} \quad (5)$$

上述关系适用于光在介质表面发生正入射的情况, 且发生于空气-样品的介质表面。具体实验中, 漫反射光谱  $R(\nu)$  是在有限的波数范围内测量得到的, 因此需要对  $\phi(\nu)$  的相位进行外推<sup>[22]</sup>:

$$i\phi(\nu) = - \frac{1}{2\pi} \int_{-\infty}^0 dt e^{-i\nu t} \int_{-\infty}^{\infty} d\nu' \ln \eta(\nu') e^{i\nu' t} + \frac{1}{2\pi} \int_0^{\infty} dt e^{i\nu t} \int_{-\infty}^{\infty} d\nu' \ln \eta(\nu') e^{i\nu' t} \quad (6)$$

通过式(1)~(5)的推导, 以及式(6)的外推方法, 最终通过式(5)~(6)可以从物质的漫反射光谱中计算吸收系数  $k$ 。

## 2.3 漫反射光谱校正方法的应用

取实验测量的 20 组葡萄糖粉末样品光谱平均值, 并进行 7 点多项式(Savitzky-Golay, SG)平滑减弱激光器功率噪声对分析结果的影响。如图 3~4 所示, 图 3 为计算后葡萄糖漫反射光谱; 图 4 为漫反射校正变换后得到的吸收光与参考数据<sup>[23]</sup>的对比。通过对比可以观察到, 实验参考所用的葡萄糖吸收光谱在  $1\ 024\ \text{cm}^{-1}$ 、 $1\ 079\ \text{cm}^{-1}$ 、 $1\ 106\ \text{cm}^{-1}$ 、 $1\ 149\ \text{cm}^{-1}$ 、 $1\ 204\ \text{cm}^{-1}$  附近出现了吸收峰, 峰值位置已在图中用虚线标出。其中 K-M 方程变换后对应吸收峰位置的透过率分别为 0.08、0.45、0.35、0.22、0.12; K-K 方程变换后对应吸收峰位置的透过率分别为 0.18、0.63、0.59、0.50、0.32。K-M 方程和 K-K 方程计算后的合成光谱在吸收峰附近也出现了峰值, 应用 O'Haver 团队开发的寻峰软件对合成光谱进行吸收峰的强度定量计算, 调整幅度、斜率、光滑宽度和自适应宽度 4 个参数约束寻峰最终的结果<sup>[24]</sup>。其中幅度设置为 1, 斜率设置为 0.005, 光滑宽度和自适应宽度设置为 10, 在光谱分析中常用峰面积表示峰强度, 其中 K-M 方程变换后对应吸收峰强度分别为 32.04、4.02、13.19、21.75、17.26; K-K 方程变换后对应吸收峰强度分别为 20.78、5.53、11.18、19.15、12.30, 峰值相关性系数分别为 0.84 和 0.91, 误差范围在  $\pm 5\ \text{cm}^{-1}$ 。这几个吸收峰的出现是由于 C-C 键和 C-O 键的拉伸以及原子群 OCH、COH、CCH 的变形劈裂成多个峰<sup>[25]</sup>。这个结果说明在不经诸如 KBr 压片法等样品预处理方法下, 通过漫反射校正方程变换, 也可以从漫反射光谱中得到样品吸收光谱的信息。

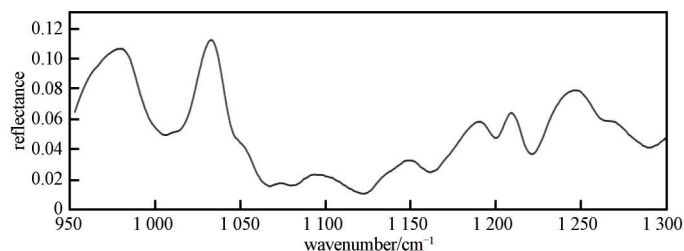


Fig.3 Measured diffuse reflectance spectroscopy of glucose

图 3 实验测量葡萄糖漫反射光谱

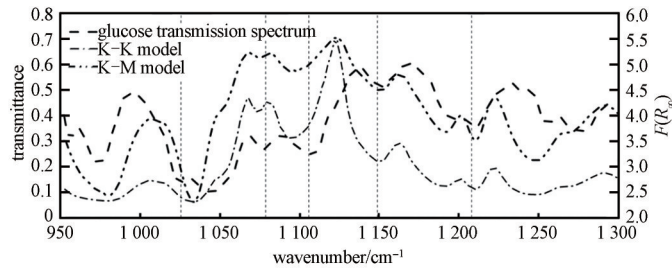


Fig.4 Corrected absorption spectra and reference spectra  
图4 漫反射校正变换后得到的吸收光谱与参考数据的对比

### 3 神经网络建模与结果分析

#### 3.1 神经网络模型选取

为识别预测混合样品中葡萄糖的质量分数和探究漫反射光谱校正方法对识别效果的影响, 本文选择了2种神经网络模型(BP神经网络和LSTM网络)建立回归任务, 对葡萄糖质量分数进行预测。

##### 3.1.1 BP神经网络和天牛须搜索算法模型

BP神经网络也称前反馈人工神经网络。该网络由输入层、隐含层和输出层构成, 其自适应过程由信号的前向传递过程和误差的反向传递过程组成, 在训练过程中循环迭代调整各层权重值, 当误差减小到预期值时输出预测结果<sup>[26]</sup>。在红外光谱数据处理中, 可以将光谱数据作为输入传递至网络模型中, 并将待输出物理参数作为网络输出, 从而构建反映光谱数据和待输出物理参数关系的回归模型。

天牛须搜索算法(Beetle Antennae Search, BAS)是在2017年提出的一种多目标函数优化算法。它的原理基于天牛的觅食天性, 因为天牛的触角对气味强度的感知很灵敏, 当天牛右边触角接收到气味强度比较强时, 就会往右飞; 反之, 天牛则会向左飞。与遗传算法、模拟退火算法、蚁群算法等类似, BAS可以在未知函数以及未知梯度的数据集上实现自动寻优, 可以显著提高寻优速度<sup>[27]</sup>。

本文将BAS和BP相结合, 在网络训练之前, 首先对初始BP神经网络设置的阈值和权重进行优化, 寻找最优的阈值和权重, 然后再应用到神经网络模型中进行训练。这样可以避免神经网络模型因为梯度下降导致陷入局部最优解的情况, 进一步提升了神经网络模型的可靠性, 也在一定程度上解决了BP神经网络容易陷入局部最小化而导致泛化能力差的问题。

##### 3.1.2 LSTM网络模型

为解决前后数据之间的依赖问题, 研究者们开发了循环神经网络(Recurrent Neural Network, RNN)用于处理时序数据, 其本质依然是BP神经网络中应用的梯度下降法。与基础的神经网络层与层连接方式不同, RNN在层与层之间的神经元处也建立了更新权重的连接。但这种方式连接的网络会导致激活函数导数发生累乘的现象, 从而导致“梯度消失”和“梯度爆炸”情况<sup>[28]</sup>。LSTM网络引入了“门”机制, 将神经网络的输入层和隐含层迁入记忆单元, 加入了遗忘机制。整个网络结构由输入门、遗忘门、输出门和循环自连接的神经元构成, 通过遗忘与输入、输出的方式有效解决了“梯度消失”和“梯度爆炸”的问题<sup>[29]</sup>。

#### 3.2 模型评价指标

为对模型预测精确度进行评价, 本文选择均方误差(Mean Square Error, MSE)、平均绝对百分比误差(Mean Absolute Percent Error, MAPE)作为误差评价指标。其中:

$$E_{MS} = \sqrt{\frac{\sum_{t=1}^n |\text{actual}(t) - \text{forecast}(t)|^2}{n}} \quad (7)$$

式中:  $n$  为目标值个数; actual为实际值; forecast为预测值。

MSE用来评判预测模型的精确度, 值越小, 模型预测结果越好。本文中, 目标值为葡萄糖质量分数, 可能会出现有部分质量分数下预测结果偏离真实值较大, 但模型整体预测评价效果好的现象, 从而导致对模型错误的评估。因此引入MAPE来评价预测结果较真实结果的平均偏离水平, 如MAPE为10, 则说明预测结果偏离真实结果10%, 其表达式如下:

$$E_{\text{MAP}} = \frac{1}{n} \sum_{t=1}^n \left| \frac{\text{actual}(t) - \text{forecast}(t)}{\text{actual}(t)} \right| \quad (8)$$

### 3.3 预测结果分析

#### 3.3.1 BAS-BP 神经网络和 LSTM 网络模型对比

本文选择的 BP 神经网络隐含层神经元节点数根据经验公式  $m = \sqrt{n+l} + \alpha (\alpha = 1, 2, \dots, 10)$  ( $n$  为输入层节点数,  $l$  为输出层节点数,  $m$  为隐含层节点数) 设置为 10 个, 初始训练权重值为 1, 训练误差设置为  $10^{-4}$ , 训练集按照 120 组混合粉末漫反射光谱总数据集 2/3 的比例随机选取, 所有数据作为最终测试集, 期望输出值为葡萄糖粉末质量分数。在选取 BAS 参数时, 由于缺少经验公式的指导, 经过多次试验调整, 最终确定初始步长  $\delta^0 = 3$ , 迭代次数  $n = 100$ 。

BP 神经网络在优化复杂目标函数时容易出现梯度平缓下降的情况, 发生这种情况时, 权值误差更新速度极小且缓慢, 大大降低了网络训练速度。考虑到方法的实用性, 为进一步提高神经网络的训练效率, 本文采用竞争性自适应重加权算法 (Competitive Adaptive Reweighted Sampling, CARS) 降低输入序列的维度。以原始光谱数据为例, 依据均方根误差 (Root Mean Square Error of Cross Validation, RMSECV) 最小值原则, 对数据集进行 50 次蒙特卡洛采样, 并进行 5 次交叉验证, 依次将提取变量加入算法中循环迭代。由图 5 可知, 循环迭代 5 次后, 在 RMSECV 达到最小值时, 样本的采样次数为 22 次, 依据此结果提取了 58 个波数下对应的数据点。图 5(a) 为蒙特卡洛采样筛选波长数量的过程; 图 5(b) 为 RMSECV 变化过程, 在曲线下降过程中消除了无关变量, 在曲线上升过程中消除了相关性强的变量; 图 5(c) 为回归系数的变化趋势, “\*” 号标记处代表筛选最佳特征波长数量的位置, 即采样次数为 22 次, 得到特征波段 58 个。

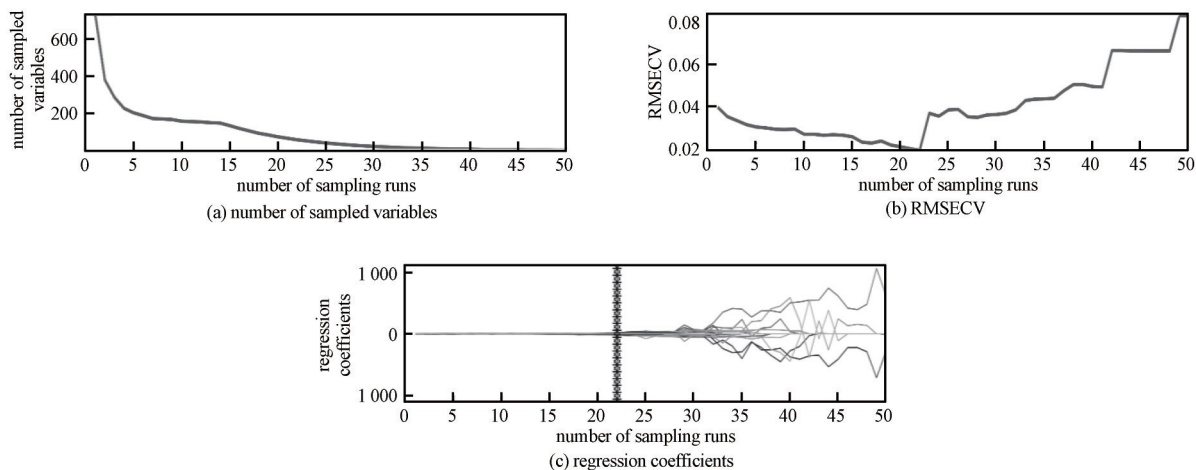


Fig.5 Wavelength variables selection process for CARS feature extraction algorithm  
图 5 CARS 特征提取算法选取波长变量过程

LSTM 网络模型的隐藏单元设置为 5 个, 为防止在训练过程中发生过拟合现象, 除去最后一层, 每一层都进行 5% 的随机权重的丢弃, 然后连接一个全连接层, 最后回归输出。训练过程使用自适应矩估计优化器, 初始学习率设置为 0.000 5。训练集按照 3.3.1 中 BP 神经网络的选取方式随机选取, 所有数据作为最终测试集, 期望输出值为葡萄糖粉末质量分数。

图 6 为原始漫反射光谱在两种网络下的预测结果, 表 1 为对原始漫反射光谱建模预测结果的评价指标。由图 6 和表 1 可知, 两种神经网络模型都实现了对目标值比较准确的预测, 但对葡萄糖质量分数小于 70% 的预测效果明显差于质量分数大于 90% 的预测效果, BAS 优化算法处理后的 BP 神经网络预测结果向较小值偏离。分析原因是由于实验选择的聚乙烯粉末粒径较葡萄糖粉末粒径小, 并且葡萄糖对红外波段的吸收比聚乙烯强, 导致聚乙烯粉末质量分数在增加的过程中, 测得的葡萄糖漫反射信号太弱, 受噪声的影响较大。LSTM 网络模型较 BP 神经网络模型 MSE 降低了 0.001, MAPE 分别降低了 5.06、21.51, 说明 LSTM 网络模型整体要优于 BP 神经网络模型。

#### 3.3.2 Kubelka-Munk 变换下 BAS-BP 神经网络和 LSTM 网络模型对比

为验证漫反射光谱校正方程对模型预测结果的影响, 首先对原始漫反射光谱进行 Kubelka-Munk 变换, 得到新的数据集, 并按照 3.3.1 中模型参数的设置方法进行建模, 得到的结果如图 7 所示。表 2 为 Kubelka-Munk 变换数据集预测结果评价指标。

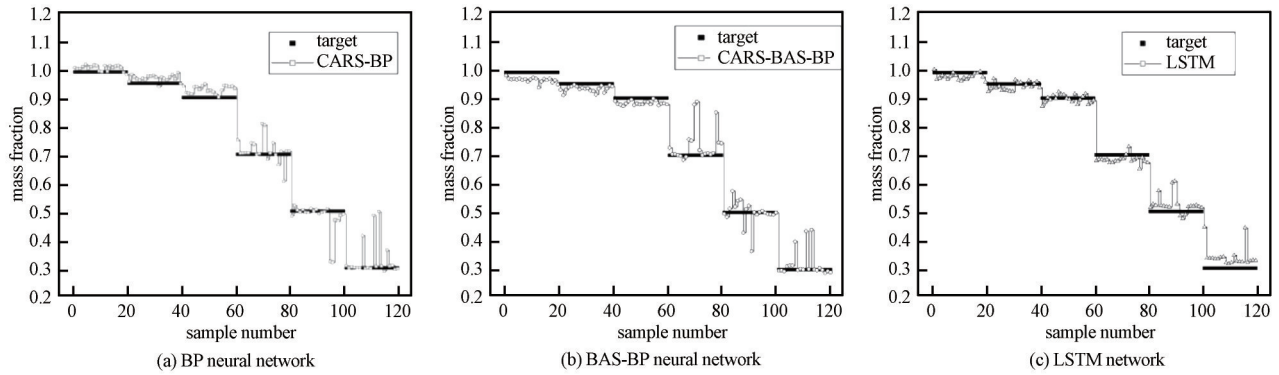


Fig.6 Predicted results for the original diffuse reflectance dataset

图6 原始漫反射数据集预测结果

表 1 原始漫反射数据集预测结果评价指标

Table1 Evaluation of prediction results of original diffuse reflection dataset

|      | CARS-BP | CARS-BAS-BP | LSTM    |
|------|---------|-------------|---------|
| MSE  | 0.000 2 | 0.002 0     | 0.001 0 |
| MAPE | 40.71   | 57.16       | 35.65   |

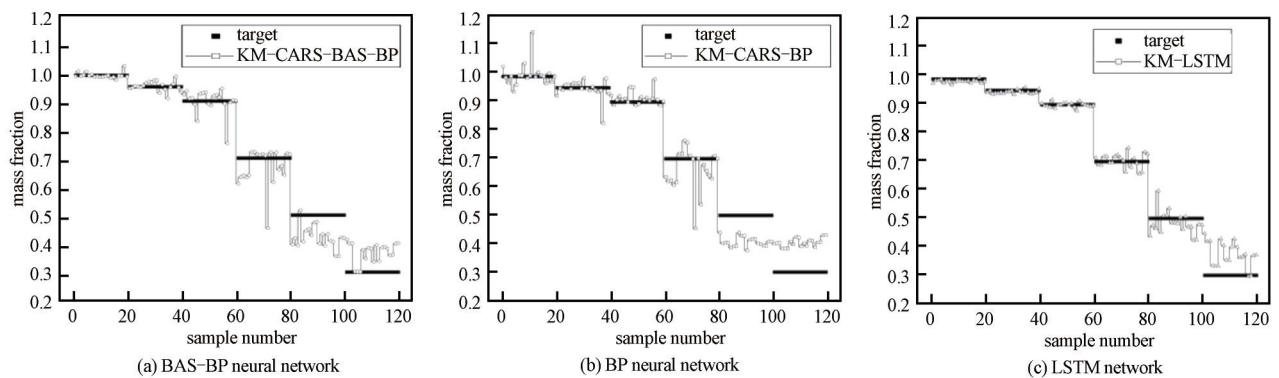


Fig.7 Predicted results from the dataset after Kubelka-Munk transform

图7 Kubelka-Munk 变换后数据集预测结果

表 2 Kubelka-Munk 变换数据集预测结果评价指标

Table2 Evaluation of prediction results after Kubelka-Munk transform

|      | KM-CARS-BP | KM-CARS-BAS-BP | KM-LSTM |
|------|------------|----------------|---------|
| MSE  | 0.005      | 0.003          | 0.002   |
| MAPE | 55.37      | 27.26          | 21.92   |

由图 7 和表 2 可知，经过 Kubelka-Munk 变换后的模型整体精确度较原始漫反射光谱模型有了提高，MAPE 在 BAS-BP 和 LSTM 模型上有了显著减小，分别减小了 29.90 和 13.7。但所有模型在葡萄糖质量分数小于 70% 的区间表现变差，分析原因是：a) Kubelka-Munk 方程一般用于混入样品浓度较低的情况，本试验中，由于葡萄糖粉末吸收强于聚乙烯粉末，在混入高浓度的聚乙烯粉末后，葡萄糖质量分数小于 70% 的区间，Kubelka-Munk 方程的适用性变差，聚乙烯样品浓度要保持在 10% 以下；b) Kubelka-Munk 方程的散射系数  $s$  依赖于样品的密度  $\rho$  以及样品粒径  $a$ ，在混合样品不均匀或者粒径大小相差较大的情况下，Kubelka-Munk 变换的效果会变差；c) 由于试验中所选葡萄糖粉末颗粒较大，用玛瑙研钵进行 2~3 min 的研磨后，使葡萄糖粉末粒径尽量接近 PE 粉末的 56~75  $\mu\text{m}$  粒径大小。若采用合适的物理研磨方法继续使样品粒径减小，Kubelka-Munk 方程变换后漫反射光谱的质量会更好<sup>[30]</sup>。

### 3.3.3 Kramers-Kronig 变换下 BAS-BP 神经网络和 LSTM 网络模型对比

为进一步验证 Kramers-Kronig 变换对模型预测结果的影响，在对原始漫反射数据集进行 Kramers-Kronig 变换之后，采用与上节相同的方式进行建模分析，得到的结果如图 8 所示。表 3 为 Kramers-Kronig 变换数据集预测结果评价指标。

由图 8 和表 3 可知，经过 Kramers-Kronig 变换后，BAS-BP 和 LSTM 网络模型的精确度有了较好的提升。BAS-BP 模型较原始的 BP 神经网络的提升主要体现在葡萄糖质量分数小于 70% 的情况下，但在葡萄糖质量分数

高于 90% 的部分, BAS-BP 模型的预测结果对比原始的 BP 神经网络模型依然向较小值偏移。LSTM 网络模型将 MSE 降到了 0.000 1, MAPE 降到了 17.08, 在所有模型中获得了最准确的预测值。分析原因是: 当激光辐照在样品表面时, 样品分子内部伸缩振动产生了吸收、反射和漫反射等复杂光学现象, 信号被激光探测器接收产生了连续带状光谱, 同时受样品粒径大小和折射率的影响, 光谱形状也会发生变化, 尤其表现在漫反射光谱中。与原子中相隔能级较大的电子跃迁产生的线状光谱不同, 分子中除了电子态之间的电子跃迁外, 还包括了能态间的振动态跃迁以及振动态内的转动态跃迁, 能级间隔较为密集, LSTM 网络可以在光谱数据中建立不同波段、峰值之间的联系关系, 从而有效地从高维度光谱数据集中反映目标物理参数的结果。对比 BP 神经网络, LSTM 网络不仅解决了前后数据之间的关系问题, 而且解决了神经网络训练中容易出现的梯度消失现象, 本试验在不增加多个隐藏层的单层网络结构下, 依然能够得到较好的预测结果, 这个特点也在最终的训练结果中得到了较好的证明。

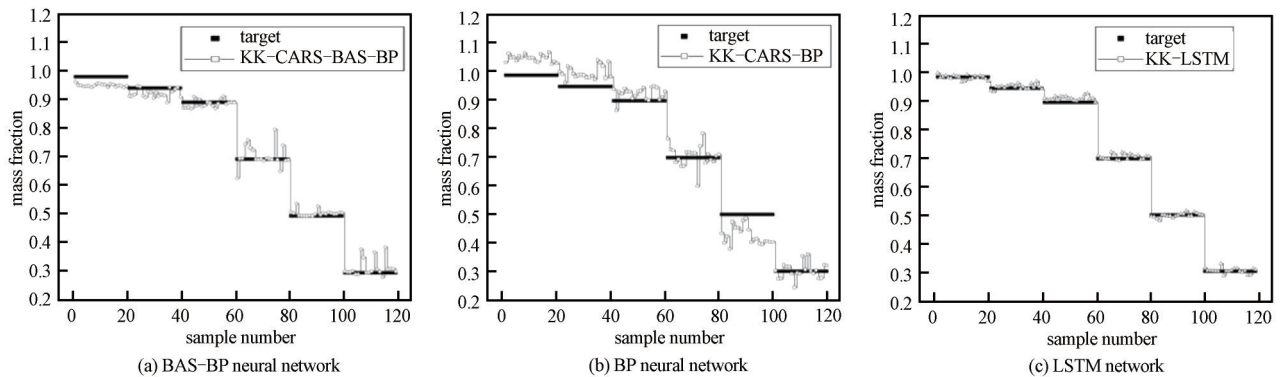


Fig.8 Predicted results for the dataset after Kramers-Kronig transform

图 8 Kramers-Kronig 变换后数据集预测结果

表 3 Kramers-Kronig 变换数据集预测结果评价指标

Table3 Evaluation of prediction results after Kramers-Kronig transform

|      | KK-CARS-BP | KK-CARS-BAS-BP | KK-LSTM |
|------|------------|----------------|---------|
| MSE  | 0.003 0    | 0.001 0        | 0.000 1 |
| MAPE | 126.04     | 61.81          | 17.08   |

## 4 结论

本文采用可调谐红外量子级联激光器作为激光源, 在漫反射光谱的实验系统下测量了葡萄糖和聚乙烯混合粉末的漫反射光谱。对实验数据应用 Kubelka-Munk 方程和 Kramers-Kronig 方程分别校正得到合成吸收光谱, 然后将数据集应用到 BP 神经网络、BAS 优化的 BP 神经网络和 LSTM 网络模型中建模分析, 对混合粉末中葡萄糖质量分数进行预测, 最终得到如下结论。

在不进行标准样品制备的情况下, 可以通过 Kubelka-Munk 方程或 Kramers-Kronig 方程变换获得葡萄糖粉末固体的合成光谱, 得到多个与吸收光谱峰值相关性较高的吸收峰, 增加了通过漫反射辐射检测物质的可靠性。通过不同网络模型对葡萄糖质量分数预测结果进行分析, LSTM 网络作为一个优秀的神经网络变种模型, 不仅继承了原始神经网络的传递特性, 且考虑了数据间的依赖性, 并解决了原始 BP 神经网络中容易出现的梯度消失问题, 在对红外光谱回归问题的分析上有较强的优势。将混合粉末原始漫反射光谱数据集作为输入时, LSTM 网络模型预测效果整体优于 BP 神经网络模型。

由于 Kubelka-Munk 方程的适用性受样品浓度、密度以及粒径大小的影响, 对原始漫反射数据集变换后预测结果在高浓度表现较差。而 Kramers-Kronig 方程通过对因果关系的建立, 所有频率下实部代表的折射率集合共同决定了虚部的消光系数, 反映了特定频率光子的吸收现象由整个频率相应的能级变化共同影响的物理意义。在本文设置的 CARS-BP、CARS-BAS-BP、LSTM、KM-CARS-BP、KM-CARS-BAS-BP、KM-LSTM、KK-CARS-BP、KK-CARS-BAS-BP、KK-LSTM 共 9 种模型中, KK-LSTM 网络模型得到了最好的预测效果。这些结论和相关机器学习技术为基于宽谱可调谐 QCL 发展未知粉末固体样品快速高光谱识别提供了依据及技术基础。

## 参考文献:

[1] 刘瑞婷. 太赫兹探测技术应用及发展研究[J]. 中国无线电, 2019(12):18-21. (LIU Ruiting. Research on application and



- development of terahertz detection technology[J]. *China Radio*, 2019(12):18–21.)
- [ 2 ] MIAO X,ZHAN H,ZHAO K,et al. Terahertz-dependent PM 2.5 monitoring and grading in the atmosphere[J]. *Scientia Sinica (Physica,Mechanica & Astronomica)*, 2018,61(10):60–69.
- [ 3 ] LI Z,ZHU Y,HAO Y,et al. Hybrid metasurface-based mid-infrared biosensor for simultaneous quantification and identification of monolayer protein[J]. *ACS Photonics*, 2019,6(2):501–509.
- [ 4 ] HENSLEY B S,DRAINE B J. Detection of PAH absorption and determination of the mid-infrared diffuse interstellar extinction curve from the sight line toward Cyg OB2–12[J]. *The Astrophysical Journal*, 2020,895(1):38.
- [ 5 ] GABRIELI F,DOOLEY K A,ZEIBEL J G,et al. Standoff mid-infrared emissive imaging spectroscopy for identification and mapping of materials in polychrome objects[J]. *Angewandte Chemie*, 2018,30(25):7655.
- [ 6 ] HOU X,LYU S,CHEN Z,et al. Applications of Fourier transform infrared spectroscopy technologies on asphalt materials[J]. *Measurement*, 2018(121):304–316.
- [ 7 ] PORTNOV A,ROSENWAKS S,BAR I. Detection of particles of explosives via backward coherent anti-Stokes Raman spectroscopy[J]. *Applied Physics Letters*, 2008,93(4):43–45.
- [ 8 ] BAWUAH P,ZEITLER J A. Advances in terahertz time-domain spectroscopy of pharmaceutical solids:a review[J]. *TrAC Trends in Analytical Chemistry*, 2021(139):116272.
- [ 9 ] CHEN L,LIAO D G,GUO X G,et al. Terahertz time-domain spectroscopy and micro-cavity components for probing samples: a review[J]. *Frontiers of Information Technology & Electronic Engineering*, 2019,20(5):591–607.
- [10] CAI X,SUSHKOV A B,SUESS R J,et al. Sensitive room-temperature terahertz detection via the photothermoelectric effect in graphene[J]. *Nature Nanotechnology*, 2014,9(10):814–819.
- [11] CHILDS D T,HOGG R A,REVIN D G,et al. Sensitivity advantage of QCL tunable-laser mid-infrared spectroscopy over FTIR spectroscopy[J]. *Applied Spectroscopy Reviews*, 2015,50(10):822–839.
- [12] 蒋涛,湛治强,沈昌乐,等. 低阈值连续波工作的 2.9 THz 量子级联激光器[J]. *太赫兹科学与电子信息学报*, 2016,14(5):657–661. (JIANG Tao,ZHAN Zhiqiang,SHEN Changle,et al. Fabrication of continuous-wave-work 2.9 THz Quantum Cascade Laser with low threshold[J]. *Journal of Terahertz Science and Electronic Information Technology*, 2016,14(5):657–661.)
- [13] DEUTSCH E R,KOTIDIS P,ZHU N,et al. Active and passive infrared spectroscopy for the detection of environmental threats[C]// *Advanced Environmental, Chemical, and Biological Sensing Technologies XI*. Baltimore, Maryland, United States: [s. n.], 2014: 91060A.
- [14] DEAN P,SHAUKAT M U,KHANNA S P,et al. Absorption-sensitive diffuse reflection imaging of concealed powders using a terahertz quantum cascade laser[J]. *Optics Express*, 2008,16(9):5997–6007.
- [15] GOLYAK I,MOROZOV A,SVETLICHNYI S,et al. Identification of chemical compounds by the reflected spectra in the range of 5.3–12.8  $\mu\text{m}$  using a tunable quantum cascade laser[J]. *Russian Journal of Physical Chemistry B*, 2019,13(4):557–564.
- [16] 张活. 基于太赫兹时域光谱技术的中药检测方法研究[D]. 西安:西安电子科技大学, 2018. (ZHANG Huo. Study on measurement methods of Chinese traditional medicine based on terahertz time-domain spectroscopy[D]. Xi'an, China: Xidian University, 2018.)
- [17] 张航,刘国海,江辉,等. 基于近红外光谱技术的乙醇固态发酵过程参数定量检测[J]. *激光与光电子学进展*, 2017,54(2):320–326. (ZHANG Hang,LIU Guohai,JIANG Hui,et al. Quantitative detection of ethanol solid-state fermentation process parameters based on near infrared spectroscopy[J]. *Laser & Optoelectronics Progress*, 2017,54(2):320–326.)
- [18] WANG L,HUI Y,JIANG K,et al. Potential of near infrared spectroscopy and pattern recognition for rapid discrimination and quantification of *Gleditsia sinensis* thorn powder with adulterants[J]. *Journal of Pharmaceutical and Biomedical Analysis*, 2018 (160):64–72.
- [19] 黄妙芬,王江颖,邢旭峰,等. 基于 LSTM 网络的海水石油污染含量遥感预测模型[J]. *广东海洋大学学报*, 2021,41(5):67–73. (HUANG Miaofen,WANG Jiangying,XING Xufeng,et al. Prediction model of petroleum pollution content in seawater based on LSTM network and remote sensing[J]. *Journal of Guangdong Ocean University*, 2021,41(5):67–73.)
- [20] ALCARAZ-DE-LA-OSA R,IPARRAGIRRE I,ORTIZ D,et al. The extended Kubelka–Munk theory and its application to spectroscopy[J]. *ChemTexts*, 2020,6(1):1–14.
- [21] FUFURIN I,TABALINA A,MOROZOV A,et al. Causality relations in analysis of diffuse reflectance spectra obtained by infrared quantum cascade laser[C]// *2019 International Conference on Optical Instruments and Technology:IRMMW-THz Technologies and Applications*. Beijing:[s.n.], 2020:114410G.
- [22] PETERSON C,KNIGHT B W. Causality calculations in the time domain:an efficient alternative to the Kramers–Kronig method[J]. *Journal of the Optical Society of America*, 1973,63(10):1238–1242.

- [23] O'HAVER T. A pragmatic introduction to signal processing with applications in scientific measurement[D]. Maryland, USA: University of Maryland at College Park, 2018.
- [24] BUSLOV D, NIKONENKO N, SUSHKO N, et al. Analysis of the results of  $\alpha$ -D-glucose Fourier transform infrared spectrum deconvolution: comparison with experimental and theoretical data[J]. Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy, 1998, 55(1): 229-238.
- [25] KALSTABAKKEN K A, HARNED A M. Spectral database for instructors: a living[J]. Journal of Chemical Education, 2013, 90(7): 941-943.
- [26] RUAN X, ZHU Y, LI J, et al. Predicting the citation counts of individual papers via a BP neural network[J]. Journal of Informetrics, 2020, 14(3): 101039.
- [27] 王甜甜, 刘强. 基于BAS-BP模型的风暴潮灾害损失预测[J]. 海洋环境科学, 2018, 37(3): 457-463. (WANG Tiantian, LIU Qiang. The assessment of storm surge disaster loss based on BAS-BP model[J]. Marine Environmental Science, 2018, 37(3): 457-463.)
- [28] DHRUV P, NASKAR S. Image classification using Convolutional Neural Network(CNN) and Recurrent Neural Network(RNN): a review[C]// The 2nd International Conference on Machine Learning and Information Processing(ICMLIP 2020). Hyderabad, India: [s.n.], 2020: 367-381.
- [29] 虞浩跃, 沈韬, 朱艳, 等. 基于双向长短期记忆网络的太赫兹光谱识别[J]. 光谱学与光谱分析, 2019, 39(12): 3737-3742. (YU Haoyue, SHEN Tao, ZHU Yan, et al. Terahertz spectral recognition based on bidirectional Long Short Term Memory network[J]. Spectroscopy and Spectral Analysis, 2019, 39(12): 3737-3742.)
- [30] 翁诗甫. 傅里叶变换红外光谱分析[M]. 北京: 化学工业出版社, 2010. (WENG Shifu. Fourier transform infrared spectroscopy analysis[M]. Beijing: Chemical Industry Press, 2010.)

#### 作者简介:

高 颂(1996-), 男, 在读硕士研究生, 主要研究方向为太赫兹与红外光谱识别检测成像. email: gaosong19960506@163.com.

阎结昀(1978-), 男, 博士, 副教授, 博士生导师, 主要研究方向为凝聚态理论.

王迎新(1981-), 男, 博士, 副研究员, 主要从事太赫兹科学与技术领域的研究工作.

解 研(1981-), 女, 博士, 工程师, 主要从事量子级联激光器应用技术研究.