

# 基于一维卷积神经网络和拉曼光谱的肺炎支原体菌株分类

赵勇<sup>1</sup>, 何梦园<sup>1</sup>, 王泊林<sup>2</sup>, 赵荣<sup>2</sup>, 孟宗<sup>1\*</sup>

1. 燕山大学电气工程学院, 河北省测试计量技术及仪器重点实验室, 河北 秦皇岛 066004

2. 燕山大学信息科学与工程学院, 河北省特种光纤与光纤传感重点实验室, 河北 秦皇岛 066004

**摘要** 肺炎支原体是造成人类呼吸系统疾病的主要原因。临床中, 患者感染不同肺炎支原体症状极为相似, 很难根据症状判别肺炎支原体类型并对症给药。因此, 准确判别肺炎支原体菌株类型对于发病机理和疾病流行病学研究以及临床精准治疗具有重要意义。拉曼光谱具有快速、高效、无污染等优点, 在生物医学领域逐渐得到越来越多研究者的关注。一维卷积神经网络(1D-CNN)是一类包含卷积运算且具有深度结构的前反馈网络, 在语音信号和振动信号分析等方面取得成功应用。提出一维卷积神经网络与拉曼光谱技术结合, 针对肺炎支原体主要基因型 M129 型和 FH 型样本的拉曼光谱数据集, 实现肺炎支原体菌株分类。利用光谱数据增强方法扩充原光谱数据集作为模型输入, 训练一维卷积神经网络模型, 解决由于小样本导致卷积神经网络数据饥渴问题; 为了得到最好的肺炎支原体分类效果并加速学习过程, 优化模型结构并确定最佳模型参数; 拉曼光谱测量时常混有高斯噪声、泊松噪声和乘性噪声, 为优化模型抗噪能力, 将原光谱分别叠加高斯噪声、泊松噪声和乘性噪声, 训练一维卷积神经网络模型并和 LDA, KNN 和 SVM 等传统算法进行比较。实验结果表明基于一维卷积神经网络方法, 对于叠加高斯噪声的光谱数据所建模型分类正确率为 98.0%, 叠加泊松噪声的光谱数据分类正确率为 97.0%, 叠加乘性噪声的光谱数据分类正确率为 97.0%, 分类正确率远高于基于一维卷积神经网络方法; 同时构造叠加 5, 15, 25, 35, 45 和 55 dBW 不同强度噪声的光谱数据集, 当噪声达到 55 dBW 时, 1D-CNN 模型仍能取得 92.5% 的分类正确率。因此, 一维卷积神经网络结合拉曼光谱技术应用于肺炎支原体菌株类型分类是可行的, 具有抗噪声能力强和分类正确率高的优点, 该研究为肺炎支原体肺炎快速诊断提供新思路。

**关键词** 肺炎支原体; 拉曼光谱; 定性分类; 一维卷积神经网络

**中图分类号:** O657.3 **文献标识码:** A **DOI:** 10.3964/j.issn.1000-0593(2022)05-1439-06

## 引言

肺炎支原体(mycoplasma pneumoniae, MP)感染是社区获得性肺炎最常见的病因之一<sup>[1]</sup>。它是引起人类呼吸道感染的常见病原体, 如支气管炎、肺炎等, 严重可诱发哮喘等疾病<sup>[2]</sup>。目前, 肺炎支原体菌株传统实验室诊断方法主要包括培养法、血清学法和聚合酶链反应(polymerase chain reaction, PCR)分析<sup>[3]</sup>。然而培养实验的培养周期长, 出结果慢; 血清学分析使用抗体作为生物标志物, 缺乏对疾病发作的敏感性<sup>[1]</sup>; PCR 技术具有较高灵敏度和特异性, 但需要熟练的操作人员、昂贵的仪器和复杂的样品预处理, 不能广泛应用于早期即时检测<sup>[4]</sup>。目前各种检测方法都存在耗时长, 对实

验条件、环境和人员要求高, 培养过程中易受细菌和真菌干扰, 诊断敏感性和特异性低等问题, 限制了其在临床上的广泛应用<sup>[5-6]</sup>。因此, 快速、灵敏、特异的肺炎支原体菌株检测方法研究具有重要临床意义。

拉曼光谱是一种非弹性散射的电磁辐射, 是分子振动和辐射之间能量交换的结果。拉曼光谱技术具有所需样品少、无需复杂预处理、不破坏样品、检测速度快且灵敏度高等特点, 广泛应用于乙肝<sup>[7]</sup>、肺癌、胃癌、肾病<sup>[8]</sup>等疾病的诊断。感染肺炎支原体的患者血液分子结构发生变化并反映在其拉曼光谱中, 为基于拉曼光谱判断肺炎支原体菌株类型提供理论依据。

随着深度学习在图像处理、语音识别等许多领域都取得成功应用, 卷积神经网络获得广泛关注和极大发展, 但目前

收稿日期: 2021-04-13, 修订日期: 2021-08-09

基金项目: 国家自然科学基金项目(62073280), 河北省自然科学基金项目(2020203010)资助

作者简介: 赵勇, 1978年生, 燕山大学电气工程学院副教授 e-mail: zhaoyong@ysu.edu.cn

\* 通讯作者 e-mail: mzyysu@ysu.edu.cn

应用于疾病诊断领域的光谱识别算法大多采用传统机器学习方法。相对于传统拉曼光谱分类算法而言,深度学习方法可以省去特征提取环节,简化光谱分类过程,提高识别准确率。一维卷积神经网络可以从包含各种特征的光谱中提取与目标分析物相关的信息。Liu 等<sup>[9]</sup>采用包括特征提取的金字塔形卷积层和用于分类的 2 个全连接层的 LetNet 变体的深度卷积神经网络对拉曼光谱数据分类的方法,在 RRUFF 矿物拉曼光谱数据库上取得很好的分类效果。Shao 等<sup>[10]</sup>使用 2 个卷积和 1 个全连接的卷积神经网络结合拉曼数据筛选前列腺癌骨转移的能力,使用五倍交叉验证方法对模型进行训练和测试,模型平均检测正确率 81.70%。李庆旭等<sup>[11]</sup>将可见-近红外透射光谱技术与 3 个卷积层和 1 个全连接层的卷积神经网络相结合,用于入孵前种鸭蛋受精信息的无损鉴别,测试集分类正确率 97.41%,高于逻辑回归、SVM 等传统方法。

本文提出将一维卷积神经网络(one-dimensional convolution neural network, 1D-CNN)模型应用到肺炎支原体菌株拉曼光谱识别问题,优化卷积核大小和数目等模型参数,针对 M129 型和 FH 型两类肺炎支原体光谱,模拟高斯噪声、泊松噪声和乘性噪声等测量拉曼光谱数据时的常见噪声,验证模型抗噪能力,通过和传统机器学习算法所建模型分类结果进行比较,证明所提出方法的有效性,为肺炎支原体肺炎快速诊断提供一个新思路。

## 1 数据集

### 1.1 肺炎支原体拉曼光谱数据

本文研究的肺炎支原体菌株拉曼光谱数据来源于 Dryad 光谱数据集(<https://doi.org/10.5061/dryad.s5h20>)<sup>[12]</sup>,具体选择两种主要的肺炎支原体菌株基因型 M129 型和 FH 型菌株的各 25 条拉曼光谱作为基础光谱样本模板。原始数据库中采集两类菌株拉曼光谱时,激光器光源功率为 28 mW,积分时间为 10 s,光谱采集范围 400~1 800  $\text{cm}^{-1}$ ,两类肺炎支原体菌株的基线校正后的平均拉曼光谱图如图 1 所示。

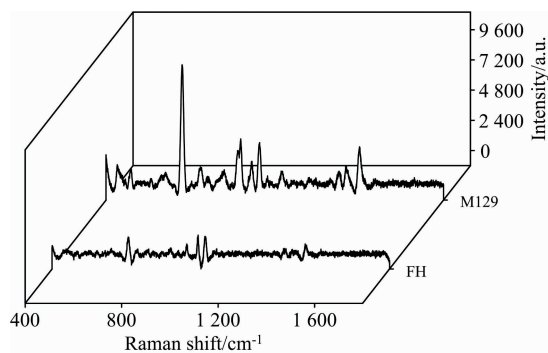


图 1 两类肺炎支原体菌株拉曼光谱图

Fig. 1 Raman spectra of two strains of mycoplasma pneumoniae

### 1.2 光谱数据增强

卷积神经网络的建立通常要求基于大规模数据库作为训练数据,充足的训练数据可以使神经网络模型充分学习到数据类别内部特征和类别间的区别,增强模型的鲁棒性,尽可

能避免过拟合现象。然而在实际应用领域中,由于临床样本的限制往往无法获得足够的拉曼光谱来训练深度学习模型。

数据增强是一种从有限的标记样本中扩大样本数量来训练神经网络,从而提高模型鲁棒性的技术。对于光谱数据,本文通过给拉曼光谱加入随机基线偏移量,设定样本光谱的不同基线斜率和随机乘性扩大光谱幅值的方法进行光谱数据增强。随机偏移量设定为样本 $\pm 0.10$ 倍光谱标准差;在 0.95~1.05 之间随机设定基线的斜率;按照样本光谱的 $1 \pm 0.10$ 倍标准差进行幅值乘性扩大。一个光谱样本进行 10 倍增强后的光谱如图 2 所示。

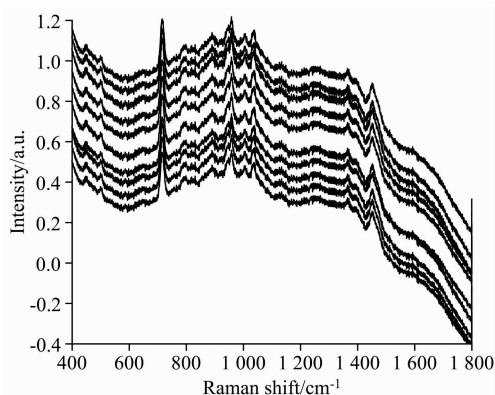


图 2 拉曼光谱数据增强

Fig. 2 Raman spectral data enhancement

### 1.3 光谱噪声叠加

实际光谱数据采集过程中,由于受传感器材料属性、工作环境、电子元器件和结构等影响,会引入各种噪声,如电阻引起的热噪声、光子噪声、暗电流噪声、光响应的非均匀性、环境噪声等。噪声以无用的信息形式出现,扰乱光谱的可观测信息,呈现与物质不相连的谱峰。传统算法很难完全去除噪声影响,因此模型的抗噪性能显得尤为重要。

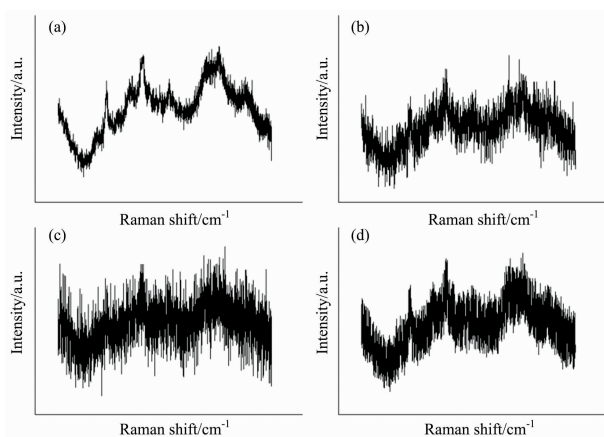


图 3 原始光谱和叠加噪声光谱

(a): 原始光谱; (b): 叠加高斯噪声光谱;

(c): 叠加泊松噪声光谱; (d): 叠加乘性噪声光谱

Fig. 3 Original spectrum and spectrum with noise

(a): Original spectrum; (b): Spectrum with Gaussian noise;

(c): Spectrum with poisson noise;

(d): Spectrum with Multiplicative noise

高斯噪声是最主要的随机噪声类型，主要是由于传感器亮度不均匀，长期工作温度过高引起；由于光具有量子特效，到达检测器表面的量子数目存在统计涨落，因此对光谱细节信息遮盖，这种由于光量子而造成的测量不确定性称为泊松噪声；乘性噪声往往由于信道不理想引起，噪声部分随信号变化而变化，且与信号是相乘关系。根据以上对噪声的分析，本文在 M129 型和 FH 型肺炎支原体菌株的原始光谱中分别叠加均值为 0，方差为 0.000 1 的高斯噪声、泊松噪声

以及方差为 0.002 的乘性噪声，叠加噪声后的光谱如图 3 所示。

## 2 模型建立与训练

根据光谱数据特点，本文构造 1D-CNN 模型进行光谱分类。整个 1D-CNN 包含 3 个卷积层和 3 个池化层，其网络结构如图 4 所示。

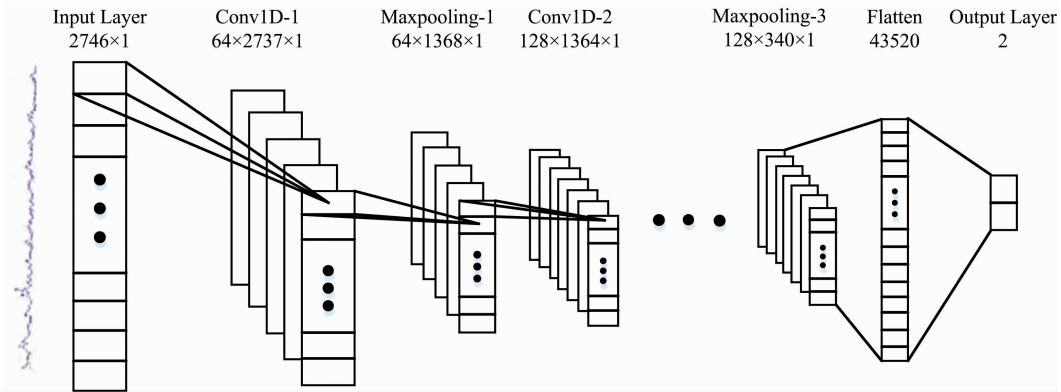


图 4 一维卷积神经网络拉曼光谱分类模型

Fig. 4 1-D convolutional neural network Raman spectra classification model

光谱数据以  $2746 \times 1$  的矩阵形式输入卷积层，在卷积运算层中，使用线性修正单元变体 LeakyReLU 作为激活函数，表示式如式(1)所示

$$f(x) = \begin{cases} x & x > 0 \\ \alpha x & \text{otherwise} \end{cases} \quad (1)$$

为了降低运算复杂度，卷积神经网络每层内的神经元权值共享，通过大量仿真优化比较，卷积核尺寸参数分别设定为  $10 \times 1$ 、 $5 \times 1$  和  $2 \times 1$ ，卷积核个数分别设定为 64、64 和 128。卷积可提取前一层信息的不同特征，这些不同特征共同作为下一层网络的输入数据。一维信号卷积运算公式如式(2)所示

$$y^j = f\left(\sum_i k^{ij} * x^i + b^j\right) \quad (2)$$

式(2)中， $*$  表示卷积运算， $y^j$  为第  $j$  个输出特征图， $x^i$  为第  $i$  个输入特征图， $k^{ij}$  为本层卷积运算所使用的卷积核， $b^j$  为第  $j$  个特征图的偏置。

每个卷积运算层后对应一个池化层，对卷积运算生成的特征图采样，池化层运算并未减少特征图个数，而是减小每个特征图的维度，缩减数据量，提升运算速度，本文采用最大池化(Max-pooling)对信号进行降采样，其表达式如式(3)所示。

$$f(X_k) = \max\{a_1, \dots, a_s\} \quad (3)$$

式(3)是最大池化对信号降采样计算方法，对卷积层运算所得的一个特征映射将其划分多个不重叠  $X_k$ ， $k=1, 2, 3, \dots, K$ ，区域大小为  $s$ 。

卷积及池化层获取光谱数据特征后，将其展开并输入全连接层进行分类，使用 tanh 激活函数，并在全连接层后加入应用比例为 0.5 的随机失活层，避免过拟合，加快收敛速度，

提高神经网络鲁棒性。设置批处理样本数目参数为 90，采用交叉熵损失函数，计算公式如式(4)所示

$$L(x_n, y_n) = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K y_{n,k} \ln p_{n,k} \quad (4)$$

式(4)中， $x_n$  是训练光谱数据， $y_{n,k}$  是第  $n$  个样本预测第  $k$  个数据的标签， $p_{n,k}$  是第  $n$  个样本预测第  $k$  个标签值的概率， $N$  是总共的样本数， $k$  为总标签类数。

原始光谱数据经过增强和三类噪声叠加构成建模所需的光谱数据集，将数据集随机划分为 3 部分：70% 光谱数据作为训练集，10% 光谱数据作为验证集，用于在反向传播训练过程中调整神经元权重参数；20% 光谱数据作为测试集，用于测试已训练后的网络模型性能。模型经过 200 个 epoch 训练之后的正确率和损失值曲线如图 5 所示，可以看出网络基本收敛。

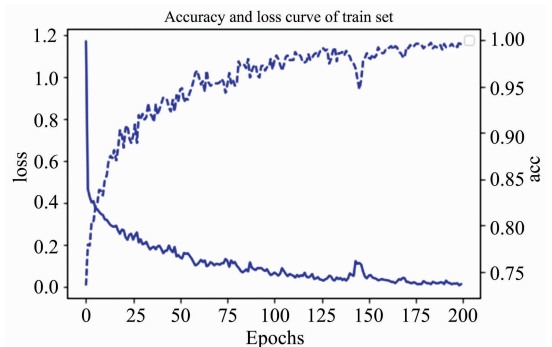


图 5 1D-CNN 模型损失率和准确率曲线

Fig. 5 Accuracy and loss curve of 1D-CNN model

### 3 模型分类结果

#### 3.1 不同算法分类结果对比

本文提出的 1D-CNN 方法与 LDA, KNN 和 SVM 三种传统方法进行比较, 结果如表 1 所示。从表 1 中可以看出,

CNN 对加入高斯噪声的 M129 和 FH 肺炎支原体菌株光谱数据所建模型分类正确率为 98%, 泊松噪声的 M129 和 FH 肺炎支原体菌株光谱数据所建模型的正确率为 97%, 乘性噪声的 M129 和 FH 肺炎支原体光谱菌株数据所建模型的正确率为 97%, 均高于传统算法所建模型的正确率。

表 1 不同算法分类结果对比

Table 1 Comparison of classification results with different algorithms

	高斯噪声			泊松噪声			乘性噪声		
	正确率	灵敏度	特异性	正确率	灵敏度	特异性	正确率	灵敏度	特异性
LDA	62.00	67.86	54.55	60.00	60.07	52.27	64.00	69.64	56.82
KNN	62.00	58.93	65.91	68.00	62.50	75.00	66.00	51.79	84.09
SVM	65.00	58.93	72.73	69.00	62.50	77.27	67.00	60.71	75.00
CNN	98.00	96.43	100.00	97.00	98.21	95.45	97.00	100.00	93.18

为了进一步比较 1D-CNN 模型和传统算法对不同种类噪声的抗噪能力, 分别得到 KNN, SVM, LDA 和 CNN 模型的混淆矩阵, 如图 6 所示。从图 6 中可以看出, 基于 1D-CNN 方法所建模型的误判个数最少, 同时 1D-CNN 模型对肺炎支原体菌株类型的灵敏度和特异性均高于其他算法所建模型, 结果表明 1D-CNN 模型相比传统算法模型在抗噪方

面具有明显优势。

#### 3.2 不同分类算法 ROC 曲线对比

对于肺炎支原体菌株类型的定性分析, ROC 曲线下面积 AUC(area under the curve)与准确率呈正相关, AUC 值越大, 模型准确率越高。不同算法针对不同种类噪声的 ROC 曲线对比结果如图 7 所示。由图 7(a)可知, 加入高斯噪声后,

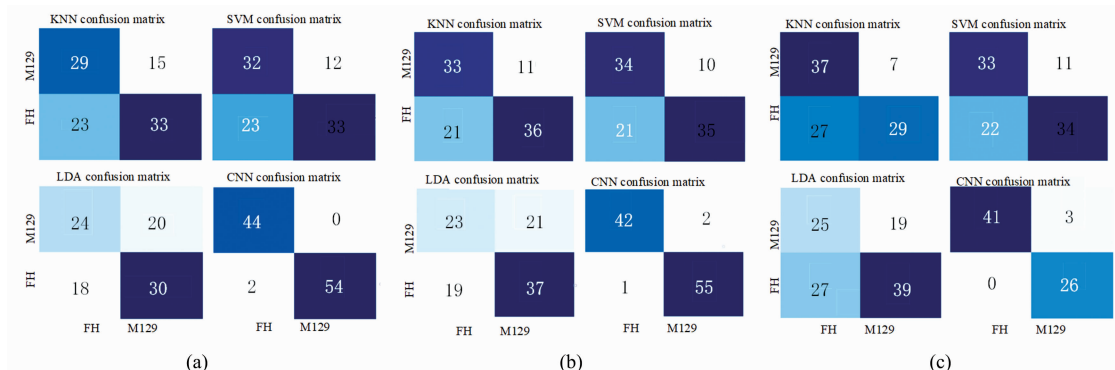


图 6 不同算法针对三种噪声的混淆矩阵对比

(a): 高斯噪声; (b): 泊松噪声; (c): 乘性噪声

Fig. 6 Confusion matrix of different algorithms with three kinds of noises

(a): Gaussian noise; (b): Poisson noise; (c): Multiplicative noise

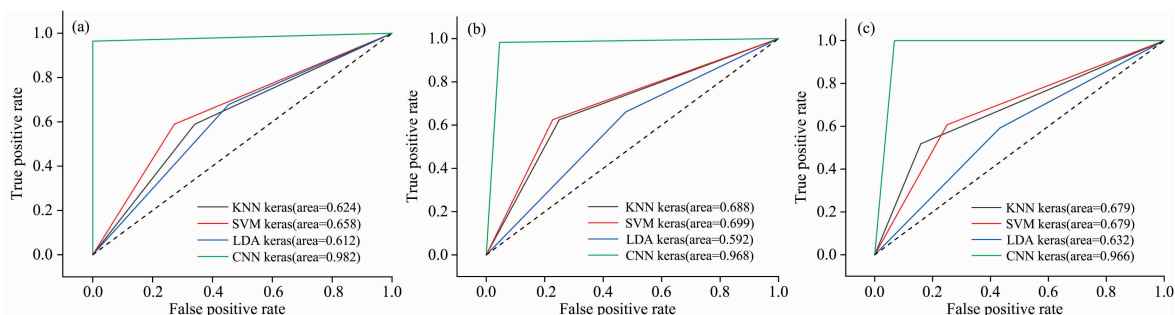


图 7 叠加三类噪声的不同算法 ROC 曲线对比

(a): 高斯噪声; (b): 泊松噪声; (c): 乘性噪声

Fig. 7 ROC curves of different algorithms with three kinds of noises

(a): Gaussian noise; (b): Poisson noise; (c): Multiplicative noise

CNN模型的AUC值为0.982, 传统分类算法LDA, KNN和SVM模型的AUC值分别为0.612, 0.624和0.658; 由图7(b)和(c)可知, 加入泊松噪声和乘性噪声后, CNN模型的AUC值也均高于其他算法。

### 3.3 叠加不同强度噪声结果对比

噪声强度对拉曼光谱定性分析模型提出更高要求, 因此

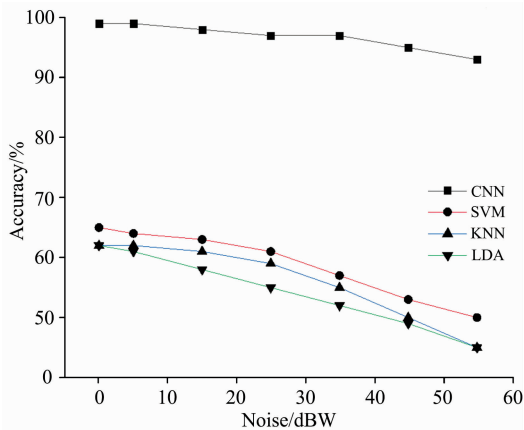


图8 不同算法针对不同强度噪声分类结果

Fig. 8 Classification results with different algorithms for different intensity of noise

在肺炎支原体菌株原始拉曼光谱中分别叠加5, 15, 25, 35, 45和55 dBW六种不同强度的高斯白噪声, 以测试模型的抗噪声性能, 四种算法在不同噪声强度下的分类结果对比如图8所示。

从图8可以看出, 随着噪声强度的不断增大, LDA, KNN和SVM算法所建模型分类正确率下降幅度较大, 而1D-CNN方法所建模型分类正确率变化幅度较小, 当添加噪声强度达到55 dBW时, 1D-CNN模型依然能够提取到拉曼光谱特征, 获得92.5%的分类正确率。因此, 1D-CNN所建模型抗噪性能远优于传统算法所建模型性能。

## 4 结论

为实现肺炎支原体菌株类型的准确分类, 提出1D-CNN拉曼光谱分类方法。针对小样本肺炎支原体菌株的拉曼光谱, 提出适用于拉曼光谱数据的数据增强方法, 扩充光谱数据以满足建模样本需求。同时模拟光谱采集时不同种类噪声影响, 验证模型的抗噪能力。结果表明利用拉曼光谱结合1D-CNN, 无需光谱预处理可以有效筛选信息, 同时能够更好地挖掘出光谱特征, 从而减少计算量和缩短计算时间。相比传统算法能得到更高的分类正确率, 并具有很好地抗噪能力, 具有明显的优势和重要的实际应用价值。

## References

- [1] Waites K B, Xiao L, Liu Y, et al. *Clinical Microbiology Reviews*, 2017, 30(3): 747.
- [2] He J, Liu M H, Ye Z F, et al. *Molecular Medicine Reports*, 2016, 14(5): 4030.
- [3] Parrott G L, Kinjo T, Fujita J. *Frontiers in Microbiology*, 2016, 7: 513.
- [4] Loens K, Leven M. *Frontiers in Microbiology*, 2016, 7: 448.
- [5] Miyashita N, Kawai Y, Kato T, et al. *Journal of Infection and Chemotherapy*, 2016, 22(5): 327.
- [6] Sano G, Itagaki T, Ishiwada N, et al. *Journal of Medical Microbiology*, 2016, 65(10): 1105.
- [7] Khan S, Ullah R, Khan A, et al. *Photodiagnosis and Photodynamic Therapy*, 2018, 23(9): 89.
- [8] Chen C, Yang L, Zhao J, et al. *Optik*, 2020, 203: 164043.
- [9] Liu J, Osadchy M, Ashton L, et al. *Analyst*, 2017, 142: 4067.
- [10] Shao X, Zhang H, Wang Y, et al. *Nanomedicine: Nanotechnology, Biology and Medicine*, 2020, 29: 102245.
- [11] LI Qing-xu, WANG Qiao-hua, GU Wei, et al (李庆旭, 王巧华, 顾伟, 等). *Spectroscopy and Spectral Analysis(光谱学与光谱分析)*, 2020, 40(12): 3847.
- [12] Henderson K C, Benitez A J, Ratliff A E, et al. *PLOS ONE*, 2015, 10(6): e0131831.

# Classification of Mycoplasma Pneumoniae Strains Based on One-Dimensional Convolutional Neural Network and Raman Spectroscopy

ZHAO Yong<sup>1</sup>, HE Men-yuan<sup>1</sup>, WANG Bo-lin<sup>2</sup>, ZHAO Rong<sup>2</sup>, MENG Zong<sup>1\*</sup>

1. School of Electrical Engineering, Yanshan University, The Key Laboratory of Measurement Technology and Instruments of Hebei Province, Qinhuangdao 066004, China
2. School of Information Science and Engineering, Yanshan University, The Key Laboratory for Special Fiber and Fiber Sensor of Hebei Province, Qinhuangdao 066004, China

**Abstract** Mycoplasma pneumoniae is the main cause of human respiratory diseases. Clinically, the symptoms of patients infected with different mycoplasma pneumoniae are very similar, so it is difficult to distinguish the type of mycoplasma pneumoniae according to the symptoms and give medication. Therefore, the accurate identification of mycoplasma pneumoniae strain type is of great significance for the pathogenesis and epidemiological research of the disease, and accurate clinical treatment. Raman spectrum has been paid more and more attention because of its advantages of fast-speed, high efficiency, pollution-free and non-destructive analysis. One-dimensional Convolution Neural Network (1D-CNN) is a kind of pre-feedback network with a deep structure, including Convolution operation. It has been successfully applied in the analysis of speech and vibration signals. The combination of the One-Dimensional Convolution Neural Network and the Raman spectral data of the main genotypes of mycoplasma pneumoniae M129 and FH were used as the research objects to realize mycoplasma classification pneumoniae strains. The spectral data enhancement method expands the original spectral data set, and the one-dimensional convolution neural network model was trained, and the problem of data hunger of convolutional neural network caused by small samples was solved. In order to obtain the best classification effect of mycoplasma pneumoniae and accelerate the learning process, the model structure was optimized, and the best model parameters were determined. Gaussian noise, Poisson noise and multiplicative noise are often mixed in Raman spectral measurement. Gaussian noise, Poisson noise and multiplicative noise are often mixed in Raman spectral measurement. In order to optimize the anti-noise ability of the model, Gaussian noise, Poisson noise and multiplicative noise were superimposed on the original spectrum respectively, and the 1D-CNN model was trained and compared with the models built by traditional algorithms such as LDA, KNN and SVM. The experimental results show that for the Raman spectra superimposed with Gaussian noise, Poisson noise and multiplicative noise, the classification accuracy of the models based on 1D-CNN method has achieved 98.0%, 97.0% and 97.0%, respectively, which are all much higher than those of the models based on LDA, KNN and SVM algorithms. At the same time, the 1D-CNN model can achieve 92.5% classification accuracy when the noise reaches the 55 dBW interference factor, aiming at the noise with different intensities of 5, 15, 25, 35, 45 and 55 dBW. Therefore, it is feasible to apply a one-dimensional convolutional neural network combined with Raman spectrum technology to the classification of mycoplasma pneumoniae strain types, which has the advantages of strong anti-noise ability and high classification accuracy. This study provides a new idea for the rapid diagnosis of mycoplasma pneumoniae pneumonia.

**Keywords** Mycoplasma pneumoniae; Raman spectroscopy; Qualitative classification; One dimensional convolution neural network

(Received Apr. 13, 2021; accepted Aug. 9, 2021)

\* Corresponding author