

DOI: 10.19650/j.cnki.cjsi.J2514018

物理知识引导的卷积神经网络故障诊断预测方法

米洁¹, 马超¹, 周海龙¹, 甄真¹, 张健²

(1. 北京信息科技大学机电工程学院 北京 102206; 2. 北京龙科数智科技有限公司 北京 101304)

摘要:针对卷积神经网络预测滚动轴承故障中的捷径学习问题进行深入研究,并提出了一种基于物理知识引导的卷积神经网络故障诊断预测模型。采用滚动轴承数据集,对卷积神经网络在滚动轴承故障诊断模型的训练过程中出现的捷径学习问题进行了分析,揭示了捷径学习现象的存在:即使卷积网络在特定的故障数据集上达到了90%以上的精度,由于捷径学习的存在,卷积网络模型并没有学习到正确的与故障理论匹配的故障特征,而是学习到了错误的特征频率或频谱图中的波形形态。对故障诊断中捷径学习现象的产生机制进行了分析,揭示捷径学习的产生机制和决策规则:卷积神经网络的捷径学习行为主要源于数据集中由背景噪声、装配等因素导致的捷径机会,模型倾向学习简单特征组合,以及综合误差导致的数据统计偏差。由于故障数据集本身无法对深度神经网络模型的学习产生足够的约束,基于滚动轴承特征频率,设计了基于轴承故障特征的敏感频带,通过带通滤波器生成物理知识数据,构建物理引导信息,输入卷积神经网络模型,引导模型学习正确的故障特征。经实验验证,基于物理知识引导的卷积神经网络能够有效避免捷径学习问题,准确提取故障核心特征,提高故障诊断和预测准确度,提升了卷积网络故障诊断模型的可信度,在航空航天等领域高端装备故障诊断中具有应用前景。

关键词: 卷积网络; 捷径学习; 物理知识引导; 故障诊断; 故障特征频率

中图分类号: TH165 TH17 **文献标识码:** A **国家标准学科分类代码:** 520.20

Physical-guided convolutional neural network model for fault diagnosis

Mi Jie¹, Ma Chao¹, Zhou Hailong¹, Zhen Zhen¹, Zhang Jian²

(1. College of Mechanical and Electrical Engineering, Beijing Information Science and Technology University, Beijing 102206, China; 2. Beijing Longke Intelligence Technology Co., Ltd., Beijing 101304, China)

Abstract: This paper conducts an in-depth study on the problem of shortcut learning in convolutional neural networks for predicting rolling bearing faults, and proposes a physics-guided convolutional neural network model for fault diagnosis and prediction. Using rolling bearing datasets, this study analyzes the shortcut learning problem that occurs during the training of CNN-based rolling bearing fault diagnosis models, and reveals the existence of the shortcut learning phenomenon: even though the convolutional network achieves an accuracy of over 90% on a specific fault dataset, due to the presence of shortcut learning, the model fails to learn the correct fault features that match the fault theory. Instead, it learns incorrect characteristic frequencies or waveform patterns in the spectrogram. The study also analyzes the generation mechanism of the shortcut learning phenomenon in fault diagnosis, and reveals the generation mechanism. Shortcut learning behavior in convolutional neural networks mainly arises from shortcut opportunities in the dataset, caused by factors such as background noise and assembly, the model's tendency to learn simple feature combinations, and data statistical biases caused by comprehensive errors. Since the fault dataset itself cannot sufficiently constrain the learning of deep neural network models, this paper designs sensitive frequency bands based on bearing fault characteristics according to the characteristic frequencies of rolling bearings. It generates physics-guided data through band-pass filters, constructs physics-guided information, and inputs it into the convolutional neural network model to guide the model to learn correct fault features. Experimental verification shows that the physics-guided convolutional neural network can effectively avoid the shortcut learning problem, accurately extract core fault features, improve

the accuracy of fault diagnosis and prediction, and enhance the credibility of the convolutional network-based fault diagnosis model. It has application prospects in fault diagnosis of high-end equipment in fields such as aerospace.

Keywords: convolutional networks; shortcut learning; physical guided; fault diagnosis; fault characteristic frequency

0 引言

在航空航天、新能源、轨道交通和高端医疗装备与机器人等领域,滚动轴承是装备机械系统的核心部件,在降低摩擦与能耗、保障装备运行安全、延长装备寿命等方面发挥着关键性的作用。然而,滚动轴承在长期服役过程中受复杂环境和极端工况载荷作用,容易因发生多种故障,可导致装备性能下降甚至引发灾难性事故。因此,必须及时精准地诊断和预测滚动轴承的各类故障,这对保障高端装备安全可靠运行,提升装备效能至关重要^[1-2]。

近年来,基于深度学习的故障诊断和预测方法受到广泛关注和研究。与传统方法相比,此类方法在复杂特征提取、非线性问题处理和跨装备泛化能力等方面优势明显,尤其是能够适应复杂动态环境和极端交变工况。国内外学者围绕基于深度学习的机械系统故障诊断问题开展了大量研究工作^[3]。Feng 等^[4]提出一种基于单层随机小波卷积核的浅层特征提取方法(random wavelet, RaVEL),揭示了随机特征空间内积与信号空间的关联机制。黎国强等^[5]使用多种母小波函数构建连续小波变换知识库,采用多模态时-频特征的对比损失函数实现模型的训练,构建零样本故障模型。Tong 等^[6]建立了轻量级协调注意力卷积神经网络模型,通过特征图的可视化揭示了该网络学习的特征信号中的冲击等故障敏感特征。Guo 等^[7]将基于梯度的类激活图(gradient-weighted class activation mapping, Grad-CAM)技术用于故障诊断数据的可视化分析,以确定模型的焦点特征与先验知识一致。李学军等^[8]使用卷积网络提取故障的空间特征,用循环网络提取时序特征,融合后使用注意力机制聚焦轴承故障特征,并使用降维技术和类激活图提供对诊断结果的解释。此外,研究人员利用各种可视化技术展示故障诊断模型在输入数据中关注的位置,以提升模型的可解释性。

目前,基于卷积神经网络(convolutional neural network, CNN)的故障诊断系统在特定数据集上的识别准确率已经达到 90% 以上^[9],逐渐应用于各类装备的故障诊断^[3]。很多学者尝试通过信号处理,降低背景噪声来进一步提高模型对故障特征的识别能力。韩延等^[10]改进了经验模态分解算法,提取故障特征并进行信号重构,使用卷积网络对消除了时间序列复杂且随机的信号

的重构信号进行模型的训练与诊断。张锐等^[11]提出了一种将傅里叶-贝塞尔级数展开和基于能量的尺度空间经验小波变换相结合的齿轮振动信号降噪方法,将降噪后的信号用于训练卷积网络模型。但是,卷积网络在处理分类问题时,常常倾向于寻求捷径,而非学习预期的关键特征,导致其泛化能力不足,也降低了模型的置信度^[12-17]。这种“捷径学习”导致了双重技术瓶颈:1) 缺乏透明性的“黑盒”机制使模型在面对未见过的数据或异常情况表现不佳时,难以对其进行有效的改进和优化,导致模型调优陷入“试错式”迭代困境;2) 模型可能偏离对核心特征的学习,转而依赖数据中的其他非核心特征进行决策而导致失误,这种不透明的特征学习与决策机制限制了深度学习模型在安全性和可靠性要求极高场景中的应用。当前,对于装备机械系统故障诊断中的深度神经网络捷径学习问题仍然缺少必要的研究,成为制约深度学习方法应用的关键问题。

针对这一问题,提出基于特征频率引导的一维卷积网络(characteristic frequencies guided 1D-CNN, CFG-1D-CNN)模型。通过研究卷积网络故障诊断模型中的捷径学习行为,揭示捷径学习的产生原因与决策规则;基于轴承的特征频率,构建引导信息,引导卷积神经网络提取核心故障特征,提高故障诊断和预测精度及可靠性。

1 卷积神经网络的捷径学习分析

利用 Grad-CAM^[18]和沙普利可加性解释(SHapley additive explanations, SHAP)^[19]可解释性工具对 CNN 故障诊断模型进行分析,从微观和宏观两个角度分析并揭示 CNN 模型在故障诊断领域存在的捷径学习现象。微观角度使用 Grad-CAM 对个体样本输出进行分析,确定模型对每种类型个体故障样本关注的特征。宏观角度使用 SHAP 对数据集整体进行分析,确定模型在总体上对每种类型的故障所关注的特征。

1.1 捷径学习分析模型

基于典型一维卷积神经网络(1-dimensional convolutional neural network, 1D-CNN)故障诊断模型,分析捷径学习行为,所采用模型的结构如图 1 所示。将振动信号从时域转换成包络谱输入模型,使用 Grad-CAM 热力图以及 SHAP 摘要图观察模型对于输入信息的关注部分,揭示模型在学习过程中存在的捷径学习问题。

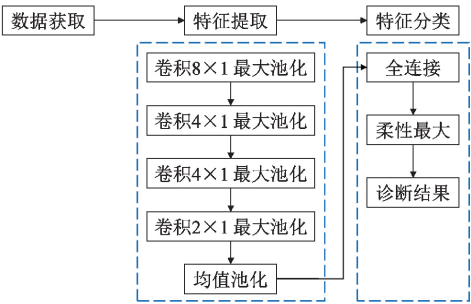


图 1 典型 1D-CNN 网络模型

Fig. 1 Typical 1D-CNN network model

1.2 捷径学习分析数据集

利用帕德博恩大学(Paderborn University, PU)轴承数据集^[20],以及凯斯西储大学(Case Western Reserve University, CWRU)滚动轴承数据中心提供的试验数据^[21]开展实验研究,包括 2 个场景:

场景 1:验证在单一工况、不同故障尺寸条件下是否存在捷径学习现象。实验采用 PU 轴承数据集中特定的样本数据;选取一级、二级损伤程度的外圈故障和内圈故障以及健康轴承测得的振动数据。按照常见的分类方法,将轴承故障数据分类为正常、内圈一级故障、内圈二级故障、外圈一级故障及外圈二级故障 5 类。根据轴承参数与转速信息,计算可以得到旋转频率和相应特征频率。

场景 2:验证在单一工况、固定故障尺寸条件下是否存在捷径学习现象。实验采用 CWRU 轴承数据集中的特定样本数据;选取故障直径为 0.177 8 mm 故障数据。样本中,轴承状态分为内圈故障、外圈故障、滚动体故障及健康轴承 4 类。根据轴承参数与转速信息,计算可以得到旋转频率和相应特征频率。

PU 和 CWRU 数据集使用的轴承型号,以及转频、内圈特征频率(ball pass frequency inner race,BPFI)、外圈特征频率(ball pass frequency outer race,BPFO)、滚动体特征频率(ball spin frequency,BSF)和保持架特征频率(fundamental train frequency,FTF)等各种特征频率如表 1 所示。

表 1 PU 和 CWRU 数据集轴承型号及特征频率
Table 1 Bearing models and characteristic frequencies in the PU and CWRU datasets

特征	PU 数据集	CWRU 数据集
轴承型号	6203	6205
转频 f_r /Hz	25	29.95
内圈特征频率 f_{BPFI} /Hz	123.68	162.19
外圈特征频率 f_{BPFO} /Hz	76.32	107.36
滚动体特征频率 f_{BSF} /Hz	49.82	67.95
保持架特征频率 f_{FTF} /Hz	9.54	11.93

1.3 PU 数据集上的捷径学习

使用 PU 数据集训练模型,按照故障类型与尺寸分类,是故障诊断面临的更具挑战性的实验场景。轴承故障尺寸的变化不会影响故障特征频率,但故障尺寸的增大会导致故障引起的冲击能量增加,振动信号中脉冲幅值与频谱特征频率及其边频带的幅值均呈显著上升趋势,这一规律在包络谱分析中尤为明显^[22]。1D-CNN 模型在 PU 数据集上的诊断精度达到 97.73%,如图 2 所示,展现出优异的分类性能。

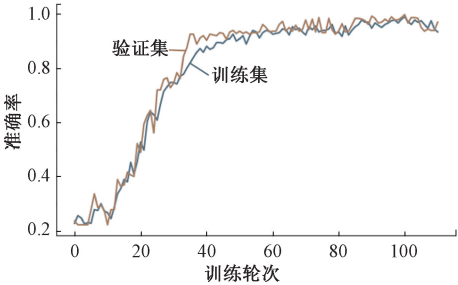
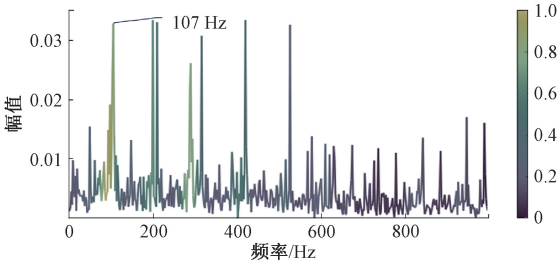


图 2 1D-CNN 模型在 PU 数据集上的精度

Fig. 2 Accuracy of 1D-CNN model on PU dataset

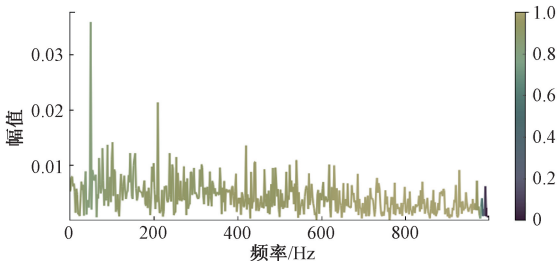
1) 个体样本中的捷径学习

使用 Grad-CAM 热力图对 PU 数据集中的样本诊断决策可视化,解析模型在输入样本中关注的关键特征区域。PU 数据集的 Grad-CAM 热力图如图 3 所示,该方法通过颜色直观量化模型对样本不同区域的注意力分布:右侧颜色条标示 0~1 的归一化权重值,浅色高亮区域表征模型在单样本分析中判定的核心判别特征。针对 5 类状态的振动信号样本,展示了模型最后一个卷积层输出的特征激活热力图。每个子图对应特定类别的一个样本,通过局部特征聚焦程度揭示模型对个体样本的决策依据。



(a) 正常轴承热力图

(a) Normal bearing heatmap



(b) 内圈一级故障热力图

(b) Inner fault level 1 heatmap

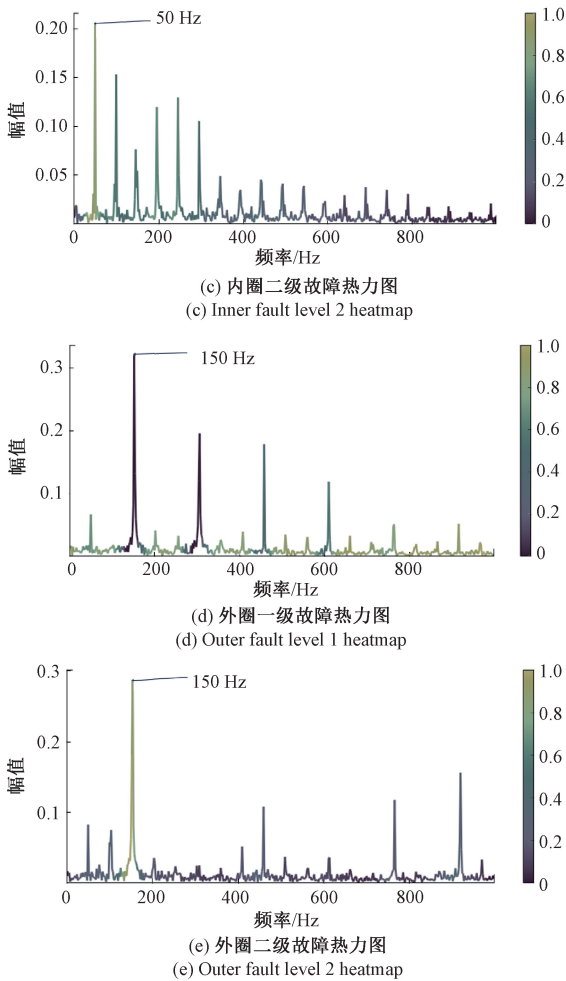


图 3 PU 数据集各类别故障热力图

Fig. 3 Heatmap for each category of the PU dataset

正常轴承热力图如图 3(a) 所示,模型未关注转频 f_r ,而是聚焦于 90~107 Hz 频段,与轴承正常、内、外圈故障均无关的频段。内圈一级和二级故障热力图如图 3(b)、(c) 所示,模型并未关注相同的特征频率,一级故障学习到包络谱的波形,而二级故障学习到滚动体特征频率 f_{BSF} 及其倍频特征。外圈一级和二级故障热力图如图 3(d)、(e) 所示,模型同样没有学习到相同的特征频率,二级故障学习到外圈特征频率 f_{BPFO} 的 2 倍频,而一级故障学习到外圈特征频率 f_{BPFO} 及其倍频除外的波形形状。

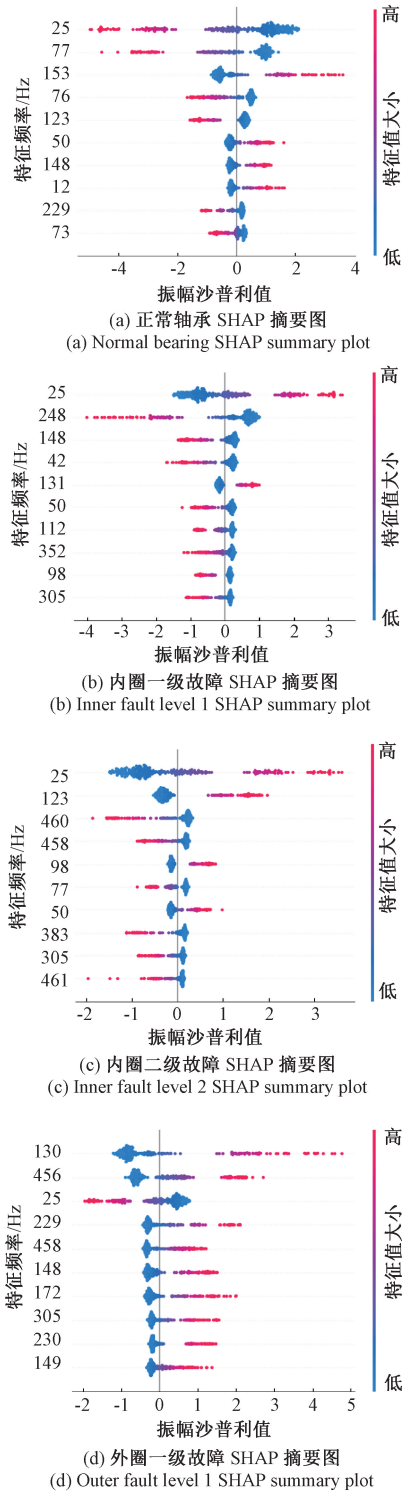
由上述热力图可见,尽管模型在数据集上表现出优异的分类性能,从个体样本上可以看到,模型的特征选择偏差揭示了 CNN 在特征学习过程中存在明显的“捷径学习”现象。

2) 数据集层面中的捷径学习

使用 SHAP 算法从各故障类别的数据集层面分析特征对模型预测的贡献度,对模型特征学习机制进行系统性评估。实验通过计算数据集中各频率分量的 SHAP

值,按照对故障分类的重要性排序。当某一故障特征频率对应的 SHAP 值显著高于其他特征时,表明该特征在模型决策过程中起到关键判别作用。这种量化分析方法不仅能够揭示模型的核心判别特征,还可通过正负向贡献值的分布情况,正值表示特征支持当前类别预测,负值表示特征抑制当前类别预测,有效观察模型决策过程。

PU 数据集各故障类别的 SHAP 摘要图如图 4 所示。



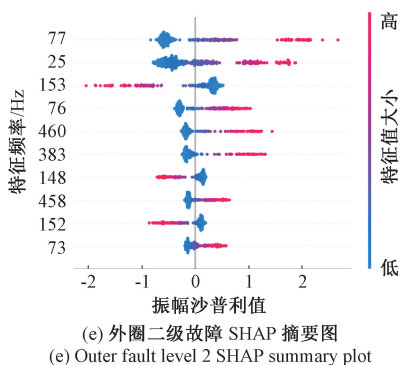


图4 PU数据集各个类别故障 SHAP 摘要图

Fig. 4 SHAP summary plots for PU dataset

图4中纵轴按全局平均 SHAP 绝对值降序排列,直观呈现特征重要性层级;横轴反映特征贡献值的大小,每个散点代表单个样本的特征值,深色表示高值,浅色表示低值。图4中展示了对模型输出影响最大的前10个特征。

正常状态轴承 SHAP 摘要图如图4(a)所示,模型错误学习到 $2f_{BPFO}$ 作为健康轴承关键特征,而 f_{BPFO} 特征频率却对该状态起到负向作用。内圈一级、二级故障 SHAP 摘要图如图4(b)、(c)所示,模型均错误的学习到 f_r 作为故障特征。此外,二级故障学习到 f_{BPFI} 特征频率;而一级故障接近 f_{BPFI} 特征频率2倍频的248 Hz 对分类起到否定作用,后续特征作用程度近似。一级故障与如图3(b)所示的 Grad-CAM 热力图相同,模型学习到包络谱的波形。外圈一级、二级故障 SHAP 摘要图如图4(d)、(e)所示,对于二级故障,模型学习到外圈故障特征频率 f_{BPFO} 以及相应高阶倍频和转频 f_r ;而一级故障没有学习到外圈故障的特征频率,且前10个特征均有较大的分类贡献,表明模型学到了包络谱的波形。

综合上述实验分析表明,卷积网络模型在没有适当约束或引导的情况下,倾向于学习“捷径”特征而非预期的故障特征。从微观和宏观两个层面分析可见,卷积网络模型存在严重的捷径学习现象。

1.4 CWRU 数据集中的捷径学习

为验证捷径学习在故障诊断中的普遍性,将1D-CNN模型在只有单一故障尺寸、相对简单的 CWRU 数据集上做进一步验证。模型在该数据集上的精度达到99.83%,展现出优异的分类性能。CWRU 数据集中内圈故障与外圈故障的 Grad-CAM 热力图如图5所示。

热力图分析结果揭示出模型在特征学习机制中同样存在显著偏差。内圈故障热力图如图5(a)所示,模型未聚焦于滚动体内圈特征频率 f_{BPFI} ,却异常关注转频 f_r 及其2、3倍频特征频率。外圈故障热力图如图5(b)

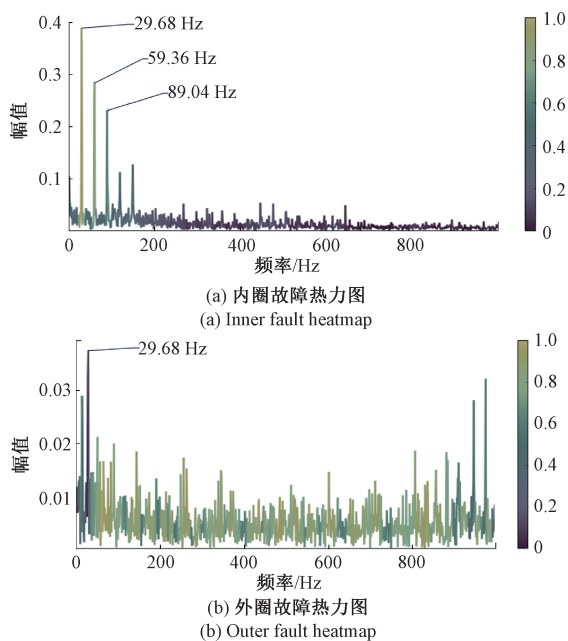
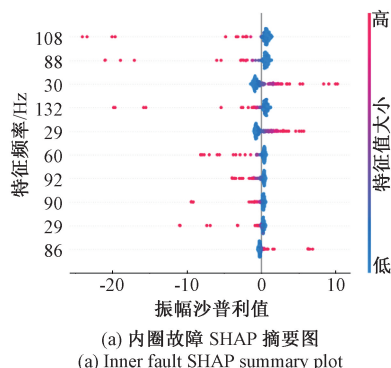


图5 CWRU数据集各类故障热力图

Fig. 5 Heatmap for each category of the CWRU dataset

所示,模型未按照预期关注外圈通过频率 f_{BPFO} 及其谐波,而是关注转频 f_r 以上的能量分布,或者说将频谱图的波形形态作为轴承外圈故障特征。对于滚动体故障和正常状态的样本,模型同样没有学习到正确的故障特征频率。

CWRU 数据集内圈故障和外圈故障的 SHAP 摘要图如图6所示,展示了对模型输出影响最大的前10个特征。内圈故障的 SHAP 摘要图如图6(a)所示,模型将转频 f_r 识别为内圈故障的重要特征,而非理论预期的 f_{BPFI} 。外圈故障的 SHAP 摘要图如图6(b)所示,虽然对模型输出影响最大的前10个特征里包含外圈故障特征频率 f_{BPFO} ,但是图6中前10个特征峰值大的点均对模型输出起到了正向影响,实际是通过捕捉连续包络谱的波形形态特征实现分类。在数据集整体上,对于滚动体故障和正常状态的样本,模型同样没有学习到正确的故障特征频率。



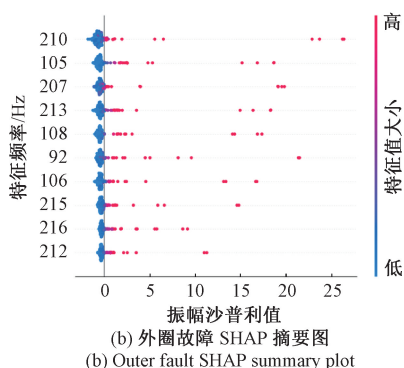


图6 CWRU数据集各类别故障 SHAP 摘要图

Fig. 6 SHAP summary plots for CWRU dataset

无论是基于数据集整体的宏观视角,还是单样本热力图的微观视角,模型在 CWRU 数据集上均表现出与理论预期不符的特征选择模式。这种“高准确率-低特征相关性”的矛盾现象表明模型存在捷径学习现象,即通过学习某些非关键特征来实现高分类性能,而非真正捕捉到研究人员期望的故障特征。这种情况说明模型的优异性能并不必然等同于正确的特征学习机制。

1.5 捷径学习类型

通过对实验的深入分析,尽管 CNN 模型能够在训练数据集上实现极高的精度性能,但并不意味着模型成功学习到正确的故障特征,或者说与故障理论相匹配的特征。当模型学习到非期望特征时,则出现了捷径学习现象。本文在故障诊断领域中针对频域(包络谱)信号的分析揭示了 2 种主要的捷径学习类型:

1) 虚假特征频率

模型未能学习到正确的特征频率,而是错误地将非故障特征频率识别为关键特征。例如,在正常状态轴承的分析中,如图 4(a) 所示,模型错误地将外圈故障特征频率的二次谐波识别为健康轴承特征。此外,在正常状态轴承的另一个分析案例中,如图 3(a) 所示,模型聚焦于与任何故障特征均无关联的 90~107 Hz 频段。这种现象表明,模型在特征学习过程中可能受到数据中某些非关键特征的干扰,从而偏离了理论预期的故障特征。

2) 频谱图中的波形

当故障对应的特征频率幅值不足以支撑分类时,如图 5(b) 所示的 CWRU 外圈故障,模型可能转向学习波形形态特征,而非具体的特征频率,这在前面的两个数据集中均有体现。这种情况下,模型并未深入理解故障特征的物理意义,而是简单地根据频谱图的波形模式进行分类。

上述两种捷径学习类型揭示了卷积网络模型在故障诊断任务中的潜在风险。即使模型在训练数据上表现出色,也不能保证其学习到了正确的故障特征。

1.6 捷径学习的产生机制与决策规则

在 CNN 模型训练过程中,捷径来源于两个方面:数据中的捷径机会,即以不同于研究人员预期的方式解决问题的可能性;以及如何将不同特征组合起来形成决策规则。本节以故障频域(轴承包络谱)特征为例,从这两个方面探讨捷径是如何出现在故障诊断领域中的。

1) 数据中的捷径机会

以 CWRU 数据集为例,CWRU 中的数据可以大致分为 3 类:(1)数据具有明显、典型的轴承故障特征;(2)数据在故障频率上显示出一定成分,但较为模糊且在频谱中不占优势;(3)数据几乎无法显现故障特征,与噪声难以区分。这种数据特征的多样性使得模型难以专注于正确的故障特征。

在影响数据质量的众多因素中,装配因素往往比故障大小或速度/负载更为关键^[21]。例如,测试轴承所在的轴与电机轴的不对中,更换故障轴承时的松紧程度差异,以及地脚螺栓的轻微松动,都可能导致尖锐的冲击脉冲产生,进而在数据中引入大量的轴转频及其谐波,如图 5(a) 所示。这些由测试环境因素产生的背景特征,为模型创造了捷径机会。对于 CNN 模型而言,熟悉的背景特征在识别过程中的重要性甚至可能超过故障本身。

此外,背景噪声也是影响数据质量的重要因素。除了机械噪声,不同传感器的噪声模式、电磁干扰,如 CWRU 数据集中自控制感应电机的变频驱动器产生的电磁干扰等,都会对数据造成影响。这些问题可能导致包络谱特征提取时出现基线偏移、特征频率附近出现噪声边带,甚至在信噪比较低时完全淹没故障特征,使模型难以学习到正确的故障特征,从而增加了捷径学习的机会。

深度学习的成功在很大程度上依赖于大量可用的标记数据。然而,与图像和自然语言等领域相比,机械故障诊断领域的标记数据相对有限。通常情况下,故障诊断模型的数据集规模较小,这使得模型更容易学习到虚假的故障特征,进一步增加了捷径学习的可能性。即使在所谓的“大数据”中,受上述因素影响而产生的系统性偏差依然存在,导致大型数据集中也可能包含许多捷径机会。由此可见,数据本身往往难以对模型产生足够的约束。

2) 决策规则的捷径

CNN 模型并不理解现实中的故障如何定义,以及用于判别的特征如何与其他特征结合。面对卷积层提取的特征,模型的判别器可以选择任何足以对给定数据集进行判别的特征组合,如图 4 所示。只要模型在训练过程中达到损失和精度要求,或训练轮数结束,则通过这些特征组合完成决策。

在典型的端到端判别学习中,这种特征组合的自由度会导致模型倾向于选择那些能够快速达到训练目标的

特征,而非真正与故障相关的特征。这种倾向使得捷径学习成为可能,因为模型可能会利用数据中的捷径特征来简化学习过程,而不是深入理解故障的物理本质。

综上所述,捷径学习的出现是数据特征和模型决策规则共同作用的结果。数据中的捷径机会为模型提供了学习非关键特征的可能性,而模型在构建决策规则时的灵活性则进一步加剧了这一现象。为了减少捷径学习的发生,需要从数据质量和模型设计两个方面入手,确保模型能够学习到真正与故障相关的特征,从而提高其在实际应用中的可靠性和鲁棒性。

2 物理知识引导的卷积神经网络模型

2.1 引入故障特征频率物理知识

深度神经网络的决策逻辑完全依赖训练数据的统计分布,而非现实世界的物理规律。在监督学习框架下,这种数据驱动的学习模式容易陷入“虚假相关陷阱”,即模型通过捕获训练数据中的统计捷径实现高精度分类。更严重的是,这种学习模式导致模型输出与物理机理的一致性缺失,制约着其在工业场景中的推广应用。

深度学习中,引入先验知识可以显著提升模型的扩展性和可解释性,同时有效避免捷径学习现象。这不仅与人类的认知能力相契合,还能打破深度学习模型的“黑盒”属性,增强其在关键应用中的可解释性和可用性。

基于轴承故障特征频率理论,设计了频谱中的特征频带滤波器组。针对实验验证使用的两个数据集,设计包含 11 个频带,分别是:

$[f_r - \Omega/2\pi, f_r + \Omega/2\pi]$, $[2f_r - \Omega/2\pi, 2f_r + \Omega/2\pi]$, $[f_{BPFO} - \Omega/2\pi, f_{BPFO} + \Omega/2\pi]$, $[2f_{BPFO} - \Omega/2\pi, 2f_{BPFO} + \Omega/2\pi]$, $[3f_{BPFO} - \Omega/2\pi, 3f_{BPFO} + \Omega/2\pi]$, $[f_{BPFI} - \Omega/2\pi, f_{BPFI} + \Omega/2\pi]$, $[2f_{BPFI} - \Omega/2\pi, 2f_{BPFI} + \Omega/2\pi]$, $[3f_{BPFI} - \Omega/2\pi, 3f_{BPFI} + \Omega/2\pi]$, $[f_{BSF} - \Omega/2\pi, f_{BSF} + \Omega/2\pi]$, $[2f_{BSF} - \Omega/2\pi, 2f_{BSF} + \Omega/2\pi]$, $[3f_{BSF} - \Omega/2\pi, 3f_{BSF} + \Omega/2\pi]$ 。

其中 $\Omega/2\pi$ 是可调阻带截止频率参数,用于控制特征带宽度。振动数据经过带通滤波器滤波,得到引入故障特征知识生成的数据。CWRU 数据集中一个外圈故障数据经过上述预处理后的故障特征引导数据如图 7 所示,展示了包络谱的振幅包络(浅色信号)和引入物理知识生成的数据(深色信号),上方的矩阵块表示 11 个敏感频段的分布,其中包含丰富的故障信息,而其他频段则与故障无关。

将引入物理知识生成的数据输入到模型中,可以有效缩小包含频段特征信息的搜索参数范围,提高搜索效率,同时实现了将先验知识融入进模型中调节学习过程,

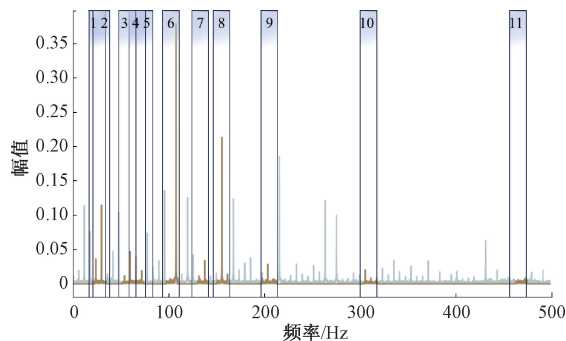


图 7 CWRU 数据的特征频带

Fig. 7 Characteristic frequency bands of CWRU dataset

与人类认知能力相契合,提升了深度学习在故障诊断场景下的可解释性。

2.2 引入物理知识的卷积网络模型

物理知识引导的智能故障诊断模型 CFG-1D-CNN, 将基于故障特征频率生成的数据与原始数据作为双通道输入进模型中进行学习,模型架构如图 8 所示。

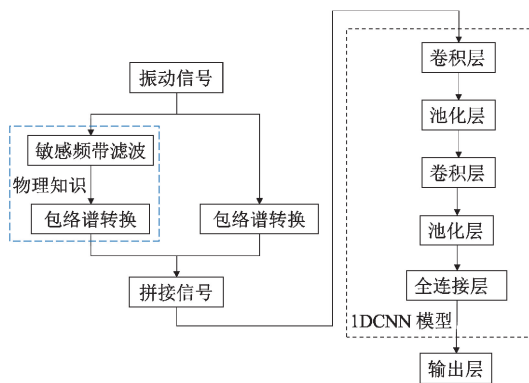


图 8 物理知识引导的模型架构

Fig. 8 Physical knowledge guided model architecture

1) 物理信息引导数据增强

物理信息引导的双通道输入设计包括两部分,第 1 部分是包络谱通道,对原始振动信号进行包络谱变换,保留完整频域信息。第 2 部分是物理引导通道,通过带通滤波提取 11 个特征敏感频段。这种双通道设计既保留了原始数据的丰富信息,又通过物理约束提供了明确的特征线索,为模型学习提供了多维度的输入表示。

2) 模型决策规则约束

正则化技术是深度学习控制模型复杂度的重要方法,通过约束参数空间,在提升模型泛化能力与抑制过拟合风险之间建立动态平衡。在 1D-CNN 模型中引入 L1 和 L2 正则化技术,以优化模型性能。

(1) L1 正则化

L1 正则化通过稀疏诱导机制生成稀疏权重矩阵,即

许多权重值会被迫变为 0。CFG-1D-CNN 模型在分类器的全连接层使用 L1 正则化实现特征选择,迫使非物理特征对应的卷积核权重趋近于 0,保留与故障特征频率相关的卷积核参数,形成稀疏特征响应模式。

(2) L2 正则化

L2 正则化通过权重衰减实现噪声抑制,防止模型过度拟合噪声数据。CFG-1D-CNN 模型(如图 9 所示)在特征提取的卷积层使用 L2 正则化实现压缩非关键特征的权重幅值,平衡特征响应的动态范围,形成平滑的特征激活模式。

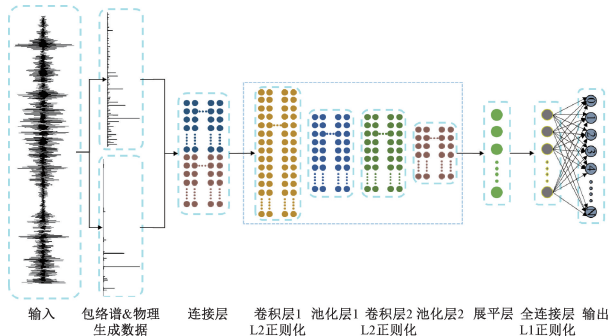


图9 CFG-1D-CNN 模型结构

Fig. 9 CFG-1D-CNN model

3 实验分析与验证

为验证所提出引入物理知识的可解释性模型的有效性和通用性,使用 3.1 节实验中的 PU 数据集和 CWRU 数据集对模型进行验证,方便对比。同时,使用实验室传动系统故障试验台采集的自有数据集进一步验证。实验使用 2.1 节设计的敏感频段对输入数据进行滤波处理,得到基于物理知识生成的数据,输入 CFG-1D-CNN 可解释性模型。

3.1 PU 数据集实验结果

CFG-1D-CNN 模型在 PU 数据集上的精度如图 10 所示,当训练轮次超过 80 轮时,模型分类精度接近 97.85%,具有良好的训练效果。

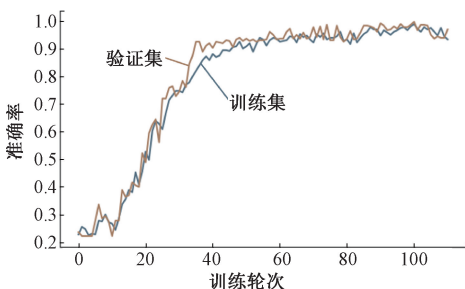
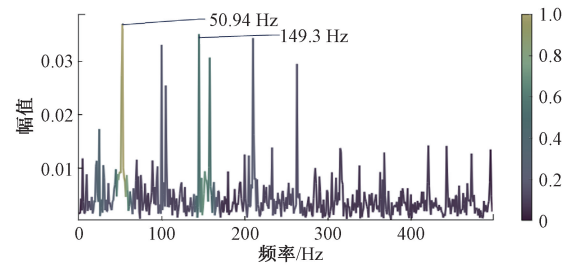


图10 CFG-1D-CNN 模型 PU 数据集的精度

Fig. 10 Accuracy of CFG-1D-CNN model on PU dataset

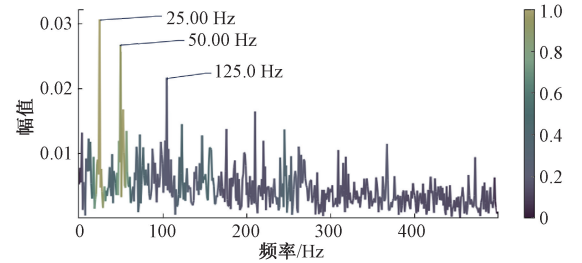
1) Grad-CAM 热力图个体样本分析

对 CFG-1D-CNN 模型中的卷积层 2 进行可视化处理,得到 5 类状态的热力图,如图 11 所示。正常轴承热力图如图 11(a)所示,模型关注转频的二次谐波 $2f_r$,即模型能够识别正常状态下的关键频率特征。内圈一级和二级故障热力图如图 11(b)、(c)所示,模型关注的特征基本一致,主要包括转频 f_r 及其谐波,以及内圈故障特征频率 f_{BPFI} 。内圈特征频率的激活强度随故障程度增加显著提升,验证了模型对故障发展阶段的幅值敏感性。外圈一级和二级故障热力图如图 11(d)、(e)所示,模型聚焦外圈的故障特征频率 f_{BPFO} 及其谐波 ($2f_{BPFO}$ 、 $3f_{BPFO}$)。



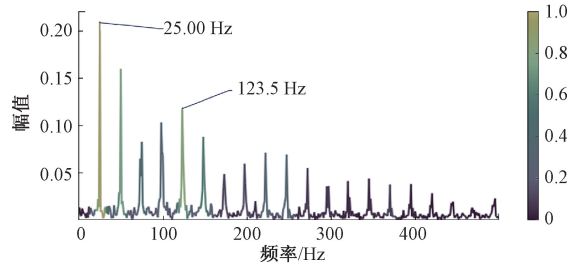
(a) 正常轴承热力图

(a) Normal bearing heatmap



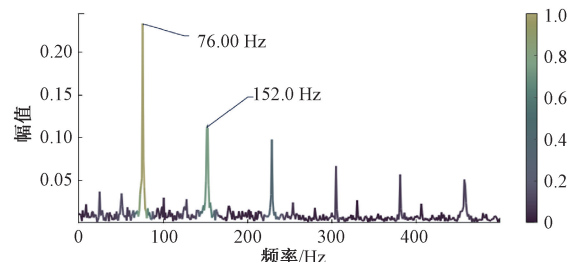
(b) 内圈一级故障热力图

(b) Inner fault level 1 heatmap



(c) 内圈二级故障热力图

(c) Inner fault level 2 heatmap



(d) 外圈一级故障热力图

(d) Outer fault level 1 heatmap

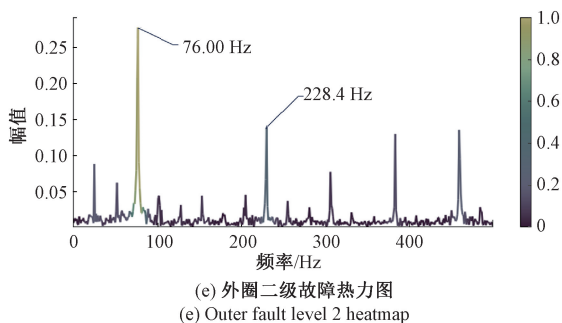


图 11 CFG-1D-CNN 模型的 PU 数据集热力图

Fig. 11 CFG-1D-CNN model's heatmap on PU dataset

综上,引入物理知识的 CFG-1D-CNN 模型能够关注与物理知识对应的故障特征频率,有效地避免捷径学习现象的发生。对于同一类型故障的不同故障程度,能够通过相同特征频率的幅值变化完成分类,符合故障机理。

2) SHAP 摘要图数据集层面分析

CFG-1D-CNN 模型在 PU 数据集上各类别的 SHAP 摘要图如图 12 所示。正常轴承的 SHAP 摘要图如图 12(a) 所示,模型主要关注转频 f_r 和 $3f_r$,是能量集中区域,幅值较低且平稳,与正常轴承振动能量集中于转频分量的物理规律相符。内圈一级和二级故障 SHAP 摘要图如图 12(b)、(c) 所示,模型主要关注转频 f_r 以及内圈故障特征频率 f_{BPFI} 。模型学习到了相同且正确的故障特征。随着故障严重程度增加, f_r 和 f_{BPFI} 的幅值对模型分类的影响显著增大。外圈一级和二级故障 SHAP 摘要图如

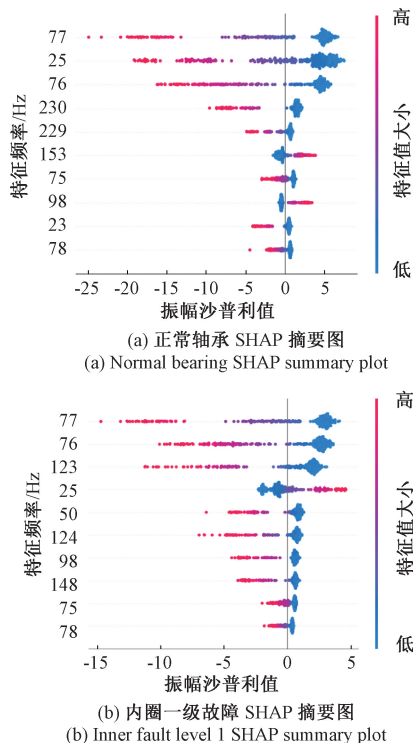
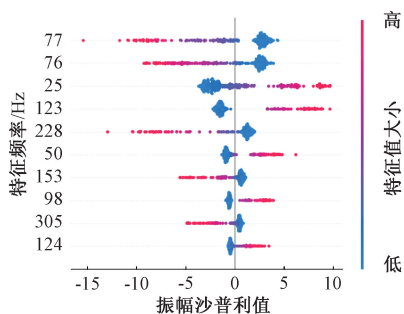
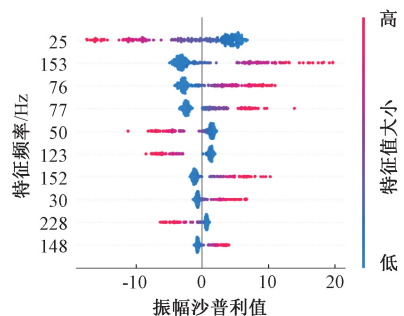
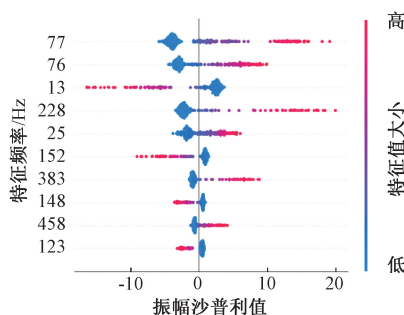
(b) 内圈一级故障 SHAP 摘要图
(b) Inner fault level 1 SHAP summary plot(c) 内圈二级故障 SHAP 摘要图
(c) Inner fault level 2 SHAP summary plot(d) 外圈一级故障 SHAP 摘要图
(d) Outer fault level 1 SHAP summary plot(e) 外圈二级故障 SHAP 摘要图
(e) Outer fault level 2 SHAP summary plot

图 12 CFG-1D-CNN 模型的 PU 数据集 SHAP 图

Fig. 12 CFG-1D-CNN model's SHAP summary plot on PU dataset

图 12(d)、(e) 所示,模型关注外圈故障特征频率 f_{BPFO} 、 $2f_{BPFO}$ 和 $3f_{BPFO}$,模型同样学习到相同且正确的故障特征。

在 PU 数据集上的实验结果表明,CFG-1D-CNN 能够学习到具有物理意义的故障特征。对不同严重程度的同一类型故障,模型能够通过相同特征频率的幅值完成分类,符合故障机理。

值得关注的是, f_{BPFO} 和 $3f_{BPFO}$ 在正常状态和内圈故障一级、二级的 SHAP 摘要图中均起到负向贡献。当这些频率的幅值增大时,模型会否定正常状态、内圈故障状态,而将其分类为外圈故障状态。同样,对于内圈故障特征频率 f_{BPFI} 也表现出类似的特点。这进一步说明引入物理知识的可解释性模型学能够学习到与物理知识对应的预期故障特征,并能够避免捷径学习现象。

3.2 CWRU 数据集实验结果

CFG-1D-CNN 模型针对 CWRU 数据集的分类精度接近 100%,展示了模型良好的训练效果。CWRU 数据集中内圈、外圈故障的 Grad-CAM 热力图如图 13 所示。

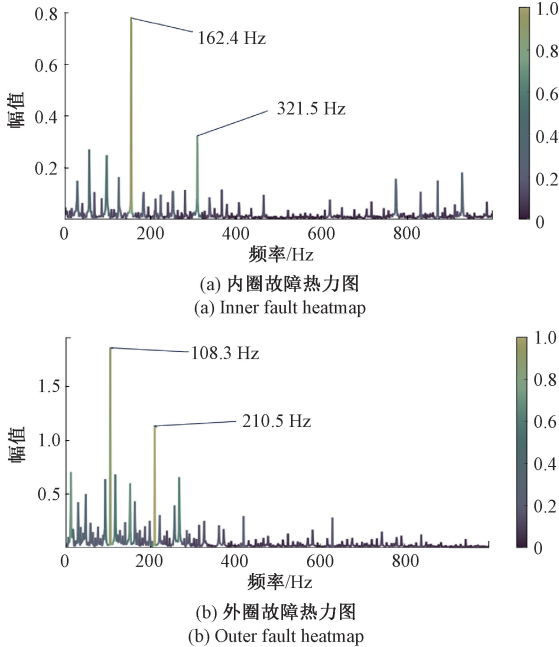


图 13 CFG-1D-CNN 模型的 CWRU 数据集热力图

Fig. 13 CFG-1D-CNN model's heatmap on CWRU dataset

内圈故障热力图如图 13(a) 所示,模型聚焦于内圈故障特征频率 f_{BPFI} 及其 2 倍频。外圈故障热力图如图 13(b) 所示,模型学习到外圈故障特征频率 f_{BPFO} 及其 2 倍频。对于滚动体故障和正常状态轴承,模型均能学习到正确的故障特征频率。改进模型在各故障类别的特征聚焦区域与理论预期完全吻合,同时也有效抑制了波形形态的捷径学习倾向。

改进模型在 CWRU 数据集上各类别的 SHAP 摘要图如图 14 所示。内圈故障的 SHAP 摘要图如图 14(a) 所示,模型学习到轴承内圈故障特征频率 f_{BPFI} 。关注的频率包括 162、160、156 和 158 Hz,这些频率的波动可能与测试中的轻微转速波动或轴承滑动有关。30 与 60 Hz 附近频率(对应 f_r 和 $2f_r$)振幅较高时对内圈故障表现正向贡献,但振幅较低时与滚动体故障成非线性关系,可能是依赖其他因素,或者与其他特征相互作用。外圈故障 SHAP 摘要图如图 14(b) 所示,模型学习到轴承外圈故障特征频率 f_{BPFO} ,关注 106 和 108 Hz。88/92 Hz(对应 $3f_r$)与内圈故障的 30/60 Hz 响应模式类似。对于滚动体故障和正常状态轴承,模型也学习到了正确的故障特征频率。

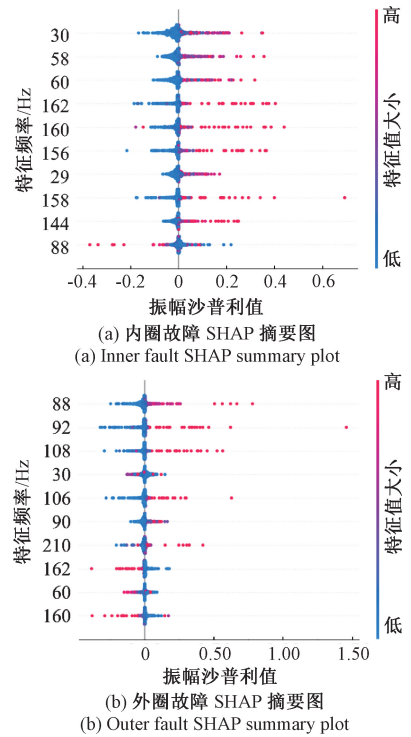


图 14 CFG-1D-CNN 模型的 CWRU 数据集 SHAP 图

Fig. 14 CFG-1D-CNN model's SHAP summary plot on CWRU dataset

3.3 传动系统故障数据集

传动系统故障数据集(transmission system fault dataset, TSFD)采集自如如图 15 所示的传动系统故障试验台,该试验台可以模拟并采集齿轮箱故障、轴承故障以及转子动平衡等故障数据。实验使用的 TSFD 数据集主要采集了正常轴承、轴承内圈故障、外圈故障,以及滚动体故障的数据。数据集的转频以及轴承相关的特征频率如表 2 所示。

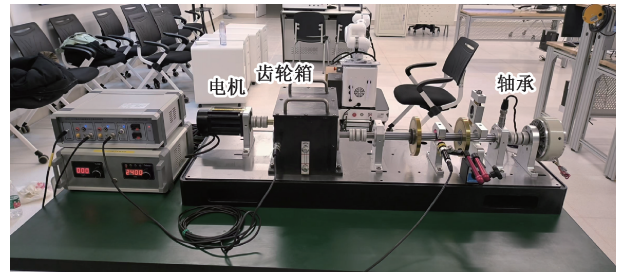


图 15 传动系统故障试验台

Fig. 15 Transmission system fault test rig

使用 1D-CNN 模型和 CFG-1D-CNN 模型对 TSFD 数据集进行分类验证,两个模型的分类精度分别为 99.34% 和 99.62%。在两个模型分类精度均接近 100% 的情况下,模型做出分类决策的依据存在着较为显著的差异。首先,观察轴承内圈故障的个体样本,如图 16 所示。1D-CNN 模型的内圈故障一个样本的热力图如图 16(a)

表 2 实验室数据集转频及轴承特征频率

Table 2 Laboratory dataset rotational frequency and bearing characteristic frequency (Hz)

特征参数	实验室轴承数据集
转频 f_r	20
内圈特征频率 f_{BPFI}	108
外圈特征频率 f_{BPFO}	72
滚动体特征频率 f_{BSF}	48
保持架特征频率 f_{FTF}	8

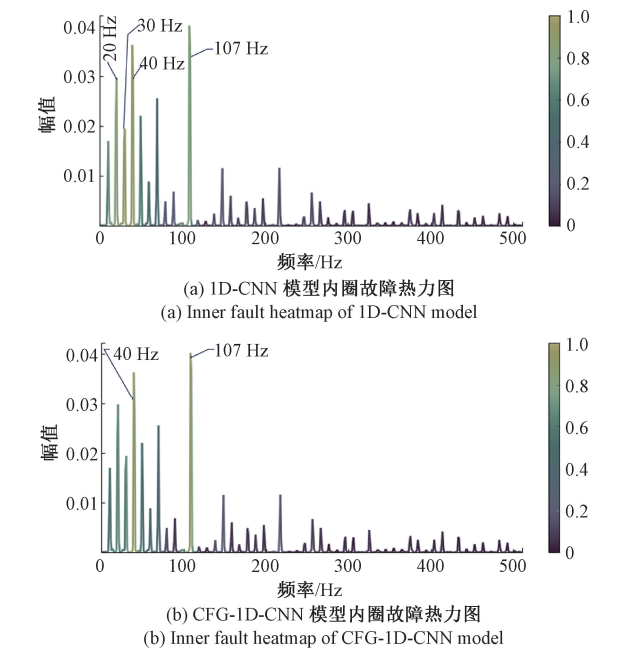


图 16 1D-CNN 和 CFG-1D-CNN 轴承内圈故障热力图对比

Fig. 16 Comparison of bearing inner ring fault heatmaps between 1D-CNN and CFG-1D-CNN models

所示,可以看到模型除了关注内圈故障特征频率 f_{BPFI} 以及转频 f_r 和二倍频 $2f_r$ 外,还关注了 30 Hz 频率,而该频率不在轴承的特征频率及相应的倍频里。CFG-1D-CNN 模型的内圈故障的热力图如图 16(b) 所示,模型关注的是内圈故障特征频率 f_{BPFI} 以及二倍转频 $2f_r$,消除了 1D-CNN 模型中存在的捷径学习现象。

1D-CNN 模型轴承内圈和外圈故障的 SHAP 摘要图如图 17 所示。从数据集层面看,对于内圈故障,如图 17(a) 所示,1D-CNN 模型除了关注轴承内圈特征频率 f_{BPFI} 和二倍转频 $2f_r$ 外,还错误的关注了外圈特征频率 f_{BPFO} 。对于外圈故障,如图 17(b) 所示,1D-CNN 模型除了关注外圈故障频率 f_{BPFO} 和二倍转频 $2f_r$ 外,还错误的关注了内圈故障特征频率 f_{BPFI} 。而 1D-CNN 模型区分内圈故障与外圈故障的一个因素是 141 Hz 特征频率,即

$2f_{BPFO}$ 。由此可见,1D-CNN 模型存在着捷径学习,即学习到错误的特征分类依据。

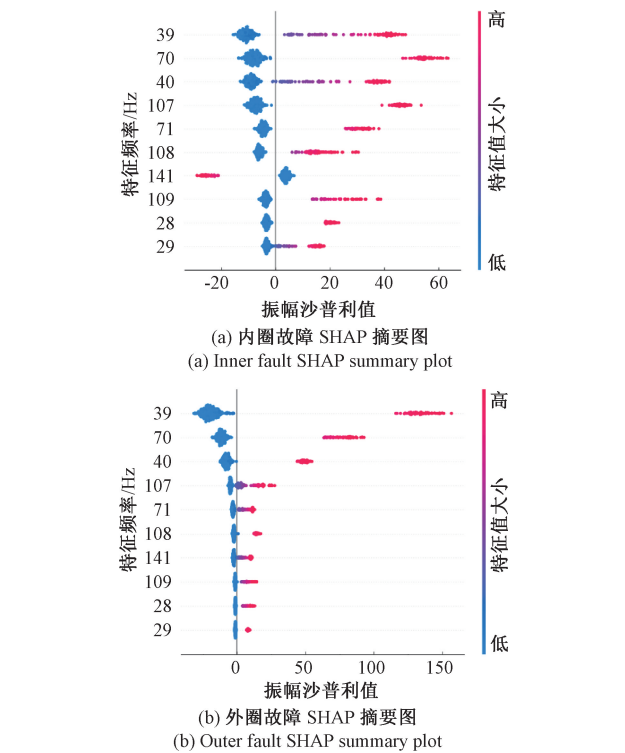
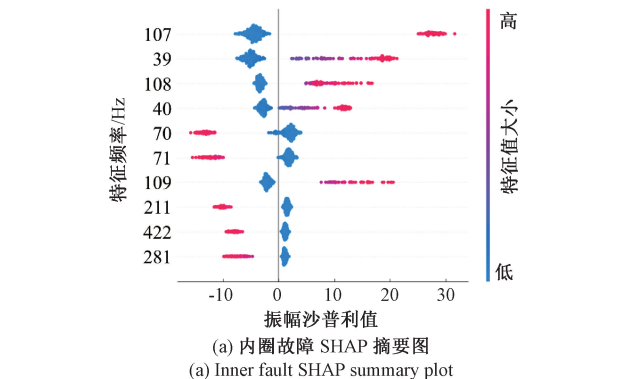


图 17 1D-CNN 模型的 TSFD 数据集 SHAP 图

Fig. 17 1D-CNN model's SHAP summary plot on TSFD dataset

CFG-1D-CNN 模型轴承内圈和外圈的 SHAP 摘要图如图 18 所示。从数据集层面看,对于内圈故障,如图 18(a) 所示,CFG-1D-CNN 模型关注内圈故障特征频率 f_{BPFI} 和二倍转频 $2f_r$,而外圈故障特征频率起到负向作用。对于外圈故障,如图 18(b) 所示 CFG-1D-CNN 模型关注外圈故障特征频率 f_{BPFO} ,而内圈故障特征频率 f_{BPFI} 起到负向作用。综合内圈故障和外圈故障 SHAP 摘要图可以观察到, f_{BPFI} 与 f_{BPFO} 是互斥的。当 f_{BPFI} 的 SHAP 值大, f_{BPFO} 的 SHAP 值小时,CFG-1D-CNN 将其分类为内圈故障,反之分类为外圈故障。这与现实的故障机理保持一致。



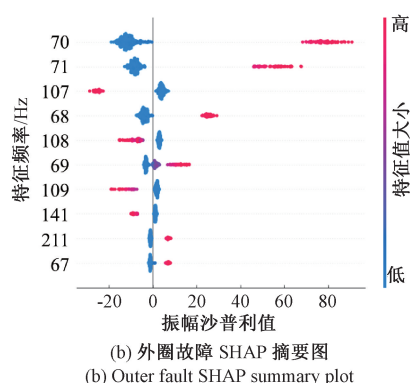


图 18 CFM-1D-CNN 模型的 TSFD 数据集 SHAP 图

Fig. 18 CFM-1D-CNN model's SHAP summary plot on TSFD dataset

虽然 1D-CNN 和 CFM-1D-CNN 模型在 TSFD 数据集上的分类精度相同,但 CFM-1D-CNN 的分类决策建立在正确的故障特征基础上,具有更好的可信度,以及应用推广的价值。

综合上面 3 个数据集的验证,通过热力图与 SHAP 分析可以看到,CFM-1D-CNN 模型成功学习到具有物理意义的特征,而没有出现捷径学习现象。这种基于物理知识的可解释性模型不仅提高了故障诊断的准确性,还增强了模型的可靠性和可解释性。

4 结 论

针对滚动轴承深度神经网络故障诊断中捷径学习问题,提出了一种基于物理知识引导的卷积神经网络模型。利用 CWRU 和 PU 两个数据集展开研究,分析了传统一维卷积神经网络在故障诊断任务中存在的基于虚假特征和非物理特征进行分类的典型捷径学习问题;揭示了这两类典型捷径学习主要源于模型倾向学习简单特征组合,以及系统噪声等因素导致的数据统计偏差。在此基础上,阐明了捷径学习的产生机制和决策规则,验证了所提出的基于物理知识引导的卷积神经网络模型能够避免捷径学习问题,准确预测故障特征。研究工作为提高深度神经网络故障诊断预测精度提供了新方法,在航空航天等领域高端装备故障诊断中具有应用价值。

参考文献

[1] 刘杰, 谭玉涛, 杨娜. 强噪声下基于 ACYCBD-MTF-MobileViT 的轴承故障诊断研究[J]. 振动与冲击, 2024, 43(24): 34-47.

LIU J, TAN Y T, YANG N. A study on bearing fault diagnosis based on ACYCBD-MTF-MobileViT under strong noise[J]. Journal of Vibration and Shock, 2024,

43(24): 34-47.

- [2] CHEN G, YUAN J L, ZHANG Y Y, et al. Enhancing reliability through interpretability: A comprehensive survey of interpretable intelligent fault diagnosis in rotating machinery [J]. IEEE Access, 2024 (12): 103348-103379.
- [3] HOANG D T, KANG H J. A survey on deep learning based bearing fault diagnosis [J]. Neurocomputing, 2019, 335: 327-335.
- [4] FENG Y, ZHENG CH Y, CHEN J L, et al. Beyond deep features: Fast random wavelet kernel convolution for weak-fault feature extraction of rotating machinery [J]. Mechanical Systems and Signal Processing, 2025, 224: 112057.
- [5] 黎国强, 魏美容, 吴德烽, 等. 零故障样本下小波知识驱动的工业机器人故障检测 [J]. 仪器仪表学报, 2024, 45(9): 166-176.
- LI G Q, WEI M R, WU D F, et al. Wavelet knowledge-driven mechanical equipment fault detection with zero-fault samples [J]. Chinese Journal of Scientific Instrument, 2024, 45(9): 166-176.
- [6] TONG J Y, LIU C, ZHENG J D, et al. Multi-sensor information fusion and coordinate attention-based fault diagnosis method and its interpretability research [J]. Engineering Applications of Artificial Intelligence, 2023, 124: 106614.
- [7] GUO L, GU X, YU Y X, et al. An analysis method for interpretability of convolutional neural network in bearing fault diagnosis [J]. IEEE Transactions on Instrumentation and Measurement, 2024, 73: 1-12.
- [8] 李学军, 刘治新, 杨同光, 等. 一种可解释性空时模型的风力发电机轴承智能诊断新框架 [J]. 仪器仪表学报, 2025, 46(2): 51-69.
- LI X J, LIU ZH X, YANG T G, et al. A new intelligent diagnosis framework for wind power insulated bearings based on spatio-temporal models of interpretable lightweight [J]. Chinese Journal of Scientific Instrument, 2025, 46(2): 51-69.
- [9] ZHANG SH, ZHANG SH B, WANG B N, et al. Deep learning algorithms for bearing fault diagnostics-a comprehensive review [J]. IEEE Access, 2020, 8: 29857-29881.
- [10] 韩延, 吴迪, 黄庆卿, 等. 基于 CNN-GraphSAGE 双分

- 支特征融合的齿轮箱故障诊断方法[J]. 电子测量与仪器学报, 2025, 39(3): 115-124.
- HAN Y, WU D, HUANG Q Q, et al. Gearbox fault diagnosis method based on CNN-GraphSAGE dual-branch feature fusion[J]. Journal of Electronic Measurement and Instrumentation, 2025, 39(3): 115-124.
- [11] 张锐,刘婷婷,王燕,等. 基于FBSE-ESEWT的齿轮故障诊断方法[J]. 电子测量与仪器学报, 2025, 39(4): 234-246.
- ZHANG R, LIU T T, WANG Y, et al. Gear fault diagnosis method based on FBSE-ESEWT[J]. Journal of Electronic Measurement and Instrumentation, 2025, 39(4): 234-246.
- [12] RUDIN C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead[J]. Nature Machine Intelligence, 2019, 1(5): 206-215.
- [13] LAPUSCHKIN S, WÄLDCHEN S, BINDER A, et al. Unmasking Clever Hans predictors and assessing what machines really learn[J]. Nature Communications, 2019, 10: 1096.
- [14] SONG R, LI Y J, SHI L D, et al. Shortcut learning in in-context learning: A survey[J]. ArXiv preprint arXiv: 2411.02018, 2024.
- [15] GEIRHOS R, JACOBSEN J H, MICHAELIS C, et al. Shortcut learning in deep neural networks[J]. Nature Machine Intelligence, 2020, 2: 665-673.
- [16] SARANRITTICHA P, MUMMADI C K, BLAIOTTA C, et al. Overcoming shortcut learning in a target domain by generalizing basic visual factors from a source domain[C]. Computer Vision-ECCV, 2022: 294-309.
- [17] KIM B, KIM H, KIM K, et al. Learning not to learn: Training deep neural networks with biased data[C]. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019: 9004-9012.
- [18] SELVARAJU R R, COGSWELL M, DAS A, et al. Grad-CAM: Visual explanations from deep networks via gradient-based localization[C]. IEEE International Conference on Computer Vision, 2017: 618-626.
- [19] LUNDBERG S M, LEE S I. A unified approach to interpreting model predictions[C]. International Conference on Neural Information Processing Systems, 2017: 4768-4777.
- [20] LESSMEIER C, ENGE-ROSENBLATT O, BAYER C, et al. Data acquisition and signal analysis from measured motor currents for defect detection in electromechanical drive systems[C]. PHM Society European Conference, 2014: 768-777.
- [21] SMITH W A, RANDALL R B. Rolling element bearing diagnostics using the Case Western Reserve University data: A benchmark study[J]. Mechanical Systems and Signal Processing, 2015, 64-65: 100-131.
- [22] 崔玲丽,张宇,巩向阳,等. 基于振动响应机理的轴承故障定量诊断及量化分析[J]. 北京工业大学学报, 2015, 41(11): 1681-1687.
- CUI L L, ZHANG Y, GONG X Y, et al. Vibration mechanism based quantitative diagnosis and quantization analysis of rolling bearing fault[J]. Journal of Beijing University of Technology, 2015, 41(11): 1681-1687.

作者简介



米洁(通信作者),1993年于华北理工大学获得学士学位,1996年于中国矿业大学(北京)获得硕士学位,2009年于北京航空航天大学获得博士学位,现为北京信息科技大学教授,主要研究方向为数字化集成设计、机械可靠性工程。

E-mail:mijie@bistu.edu.cn

Mi Jie(Corresponding author) received her B. Sc. degree from North China University of Science and Technology in 1993, received her M. Sc. degree from China University of Mining & Technology, Beijing in 1996, received her Ph. D. degree from Beihang University in 2009. Now she is a professor at Beijing Information Science and Technology University. Her main research interests include digital integrated design and mechanical reliability engineering.



马超,2003年于沈阳航空工业学院获学士学位,2009年北京理工大学获博士学位,现为北京信息科技大学副研究员;主要研究方向为转子动力学和信号检测;下肢运动动力学和运动检测。

E-mail:chaoma@bistu.edu.cn

Ma Chao received his B. Sc. degree from Shenyang Aerospace University in 2003, his Ph. D. degree from Beijing Institute of Technology in 2009. Now he is an associate researcher at Beijing Information Science and Technology University. His main research interests include rotor dynamics and signal detection.



周海龙,2021 年于集美大学获得学士学位,2025 年于北京信息科技大学获得硕士学位,现为小米科技有限责任公司工程师,主要研究方向为三维结构设计。

E-mail:2319736435@ qq. com

Zhou Hailong received his B. Sc. degree from Jimei University in 2021, his M. Sc. degree from Beijing Information Science and Technology University in 2025. Now he is an engineer at Xiaomi Corporation. His main research interest is three-dimension structure design.



甄真,2013 年于北京信息科技大学获得学士学位,2016 年于北京信息科技大学获得硕士学位,现为北京信息科技大学助理研究员,主要研究方向为机械可靠性工程。

E-mail:15810912087@ 163. com

Zhen Zhen received her B. Sc. degree from Beijing Information

Science and Technology University in 2013, her M. Sc. degree from Beijing Information Science and Technology University in 2016. Now she is an assistant researcher at Beijing Information Science and Technology University. Her main research interest is mechanical reliability engineering.



张健,1993 年于河北理工大学获得学士学位,2000 年于中国矿业大学北京研究生院获得硕士学位,现为北京龙科数智科技有限公司顾问,主要研究方向为故障预测与健康

管理。

E-mail:jacob71_zhang@ 163. com

Zhang Jian received his B. Sc. degree from North China University of Science and Technology in 1993, his M. Sc. degree from China University of Mining & Technology in 2000, Beijing. Now he is a consultant in Beijing Longke Digital Intelligence Technology Co., Ltd. His main research interests include prognostics and health management.