

DOI:10.3969/j.issn.1005-202X.2026.03.016

医学生物信息

基于图神经网络和傅里叶变换的分子属性预测

刘昱,任真

甘肃中医药大学医学信息工程学院,甘肃兰州730000

【摘要】当前的分子属性模型集中在预训练任务的改进,未对适用于分子属性预测的基础模型进行探索。本研究将无权图和有权图分别放入基于谱域和基于空域的图卷积神经网络,并使用傅里叶滤波器对节点特征进行去噪,最终模型在MoleculeNet的相关数据集上进行验证和消融实验。结果表明,该模型优于所有未经预训练的相关模型。与经典的图神经网络相比,该模型在分子属性预测任务中,对于分子图的适用性更强,能更高效地捕捉分子的结构和属性信息。

【关键词】分子属性预测;图神经网络;傅里叶变换

【中图分类号】R318;R34

【文献标志码】A

【文章编号】1005-202X(2026)03-0381-05

Molecular property prediction based on graph neural network and Fourier transform

LIU Yu, REN Zhen

School of Medical Information Engineering, Gansu University of Chinese Medicine, Lanzhou 730000, China

Abstract: Current molecular property models primarily focus on optimizing pre-training tasks, yet they lack exploration of foundational models suitable for molecular property prediction. Therefore, a novel model that integrates unweighted and weighted graphs into spectral-domain and spatial-domain graph convolutional neural networks is proposed, with Fourier filters employed to denoise node features. The proposed model is validated on relevant datasets from MoleculeNet, with ablation experiments conducted to confirm its effectiveness. Experimental results demonstrate that the model outperforms all non-pretrained baseline models. Compared with classical graph neural networks, the proposed model exhibits enhanced adaptability to molecular graphs and enable more efficient capture of molecular structural and property information, making it a promising approach for molecular property prediction tasks.

Keywords: molecular property prediction; graph neural network; Fourier transform

前言

近年来,基于机器学习的药物发现领域取得显著进展,其中分子属性预测作为关键环节,吸引了众多研究者的关注^[1-3]。早期研究多将化学分子建模为简化分子线性输入系统(Simplified Molecular Input Line Entry System, SMILES)字符串,并借助语言模型进行分子属性预测,然而,该方法存在明显局限性,难以直观呈现分子中各原子之间的化学键关系^[1,4]。鉴于此,分子图逐渐成为分子属性预测的主要工具。分子图作为一种特殊的图结构,与图神经

网络具有天然的契合性。图神经网络以图作为输入,根据卷积核的不同,通常可分为谱域类神经网络和空域类图神经网络两大类。谱域类图神经网络主要通过图傅里叶变换对图拉普拉斯矩阵进行复杂变换,即使在未经训练的情况下,也能展现出一定的性能^[5];相比之下,空域类图神经网络则侧重于通过邻域内节点的物理度量或节点向量的注意力分数进行聚合,但其性能的发挥通常依赖于良好的训练过程^[6]。图卷积网络(Graph Convolutional Network, GCN)在分子属性预测中得到广泛应用,但其受限于1-WL(Weisfeiler-Lehman)的表达能力,且深层神经网络模型并不适用于分子属性预测任务^[7-8]。傅里叶变换作为一种强大的工具,能提供频域视角,其傅里叶滤波器可以将频域中的特征与原始空域特征进行有效融合,从而形成更具判别力的特征表示。将频域中处理后的特征与原始特征相结合不仅可以增强模型对图数据的理解和表达能力,还能在保留原始

【收稿日期】2025-04-02

【基金项目】甘肃省自然科学基金(23JRRA1719)

【作者简介】刘昱,硕士研究生,研究方向:图神经网络、分子表征学习,
E-mail: 1587221637@qq.com

【通信作者】任真,副教授,硕士生导师,研究方向:医学信息处理、深度学习, E-mail: rz@gszy.edu.cn

特征局部信息的基础上,引入频域特征的全局视角,为分子属性预测提供更丰富的信息^[9]。

现有方法在处理复杂分子结构时仍面临挑战,如何更有效地表示分子结构并提高分子属性预测的准确性,仍是当前药物发现领域待解决的问题。本研究探索两种不同卷积的图神经网络的结合方法,并将傅里叶变换应用到图神经网络,同时,模型在MoleculeNet的公开数据集^[10]上进行验证,与众多基准模型进行比较,并进行消融实验。本研究不仅为分子属性预测提供一种新的方法,还为分子属性预测模型的改进提供新思路,为药物发现领域提供一种高效的解决方案。

1 相关工作

1.1 SMILES

SMILES是一种使用文本字符串描述化学分子结构的方法,主要思路是用字母、数字和特殊符号描述分子中原子的连接方式和排列顺序。SMILES更接近自然语言的形式,因此常见的使用方式是借鉴语言模型。CHEM-BERT^[4]将SMILES送入双向编码器表示(Bidirectional Encoder Representations from Transformers, BERT)模型^[11],同时将预测分子性质(分子量、芳香环数量)作为预训练任务。但是,SMILES本身存在一定问题:SMILES强制了分子中原子的排列顺序,且未能直观表示原子之间复杂的化学键关系^[1]。

1.2 分子图

分子图的出现使得化学分子的表达更灵活和精确。分子图是指将化学分子建模为图,包括原子属性作为节点特征和原子之间的连接矩阵作为边缘特征。一般将分子图分为二维分子图和三维分子图。三维分子图将坐标纳入节点特征,可以视为点云。相比于二维分子图的边缘特征是固定的,三维分子图相关模型通常使用更复杂的球面坐标系^[12]和高阶几何特征^[13]作为边缘特征。适用于三维分子图的等变图神经网络刷新了很多分子属性预测的性能,但要求数据集有精确的三维坐标且需要高昂的计算成本。对于海量的化学分子,能直接从SMILES中提取的二维分子图无疑是更合适的选择。

本研究提到的分子图是二维分子图,二维分子图根据边缘是否具有权值可分为无权分子图和有权分子图,对于分子属性预测任务一般使用有权分子图。对于有权分子图,一般使用GCN^[5]、图注意力网络(Graph Attention Network, GAT)^[6]和图同构网络(Graph Isomorphism Network, GIN)^[7]等经典图神经网络,但研究表明深度模型更关注整体而非起决定

作用的某些局部。Wang等^[14]提出的MolCLR对分子图使用原子掩码,边删除和子图删除的数据增强策略,并通过对比学习得到更丰富的分子特征。Zhu等^[15]提出的MRL-Mol和Yang等^[16]提出的MOCO在上述基础上,添加SMILES、三维分子图和分子指纹的特征并丰富预训练任务,但性能提升有限。同时,有研究者对经典图神经网络进行分层组合,不同组合的经典图神经网络性能变化不大^[17]。综上所述,当前的基础模型对于分子图的表达能力是欠缺的,一个高效的适用于分子图的模型是值得深入探索的。

2 材料与模型

2.1 分子图与模型架构

本研究提出的Dual-FFTGCN模型的整体架构如图1所示,模型的输入为SMILES字符串。SMILES字符串通过RDKit库(化学信息学工具包)进行建图,分别建立无权图和有权图,无权图送入GCN通道,有权图送入FGAT通道,最终将得到的特征合并后进行预测。其中,GCN通道在每次图卷积之后都进行池化层的处理,FGAT通道只在最后一层进行显性的池化层处理。

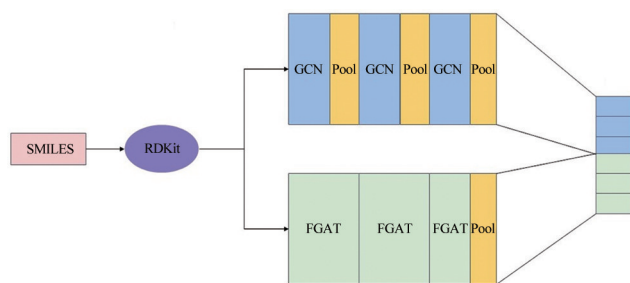


图1 模型总体架构

Figure 1 Overall model architecture

对于SMILES的二维图重建,使用表1所示的特征。对于无权图,使用表1中的节点特征,但边缘特征只使用表示连接关系的邻接矩阵,原子之间存在相应的化学键则视为存在连接;对于有权图,使用表1中的所有特征,边缘特征作为邻接矩阵的元素。事实上,RDKit可以获取原子的更多特征(电荷、杂化、芳香性和成环情况等),但是在实验中,大多数属性不适用于图神经网络的相关任务,甚至性能略有下降。

2.2 池化层

对于池化层,受傅里叶滤波器^[9]的启发,本研究开发出一种综合频域和时域的高性能自适应滤波器。对于图中的节点特征,进行实数离散傅里叶变换后提取

表1 分子图使用相关特征信息

Table 1 Relevant feature information utilized by molecular graphs

特征类型	特征名称	取值方式
节点特征	原子类型	类型不同进行独热编码
	原子序号	直接将序号作为标量
边缘特征	化学键 (单键、双键、三键和芳香键)	作为独热编码
	是否共轭	使用0、1作为布尔变量

复数信号和振幅作为特征送入门控循环单元(Gated Recurrent Unit, GRU)^[18],得到滤波系数。在频域上,使用复数信号与滤波器系数进行乘积。最终,得到恰当滤波后的节点特征信号。具体的公式如下:

$$\text{fft}_h = \text{RFFT}(h) \quad (1)$$

$$\text{fft}_h^{\text{concat}} = \text{CONCAT}(\text{REAL}(\text{fft}_h), \text{IMAG}(\text{fft}_h), \text{ABS}(\text{fft}_h)) \quad (2)$$

$$h_{\text{filter}} = \text{GRU}(\text{fft}_h^{\text{concat}}) \quad (3)$$

$$\text{fft}_h^{\text{filtered}} = h_{\text{filter}} \odot \text{fft}_h \quad (4)$$

$$h = \text{IRFFT}(\text{fft}_h^{\text{filtered}}) \quad (5)$$

其中, \odot 表示逐元素相乘。RFFT表示离散实数傅里叶变换,IRFFT表示逆离散实数傅里叶变换,GRU是一种用于序列建模的长短期记忆神经网络。在上述公式中,节点特征 h 作为输入,经过傅里叶变换转换到频域,并获取实部、虚部和振幅作为频域特征,GRU根据频域特征生成滤波器,滤波后的频域特征再通过逆变换回到时域,从而得到经过池化处理的节点特征。这一过程将频域获取的信息视为时间序列,并结合门控记忆神经网络的优势,能在图神经网络任务中对节点特征进行非线性局部特征提取,提升模型的表达能力和计算效率。

2.3 图神经网络

图神经网络根据其卷积核的不同可分为基于谱域的图卷积神经网络和基于空域的图卷积神经网络。使用经典的GCN作为谱域图卷积网络,而经典的空域图卷积神经网络GAT具有一定的局限性。GAT没有对于边缘特征的相应处理且过度依赖多注

意力的实现。本研究从EGNN^[19]中得到启发,FGAT同时使用傅里叶变换提取两个节点属性之间的权重并以恰当的方式更新节点特征。具体的公式如下:

$$h^{\text{in}} = \text{POOL}(h) \quad (6)$$

$$m_{ij} = \phi_c(h_i^{\text{in}}, h_j^{\text{in}}, e_{ij}) \quad (7)$$

$$m_i = \sum_{j \in N(i)} m_{ij} \quad (8)$$

$$h_i = \phi_h(h_i, m_i) \quad (9)$$

其中,POOL为上一节所定义的池化层, e_{ij} 是节点对之间的边缘特征, ϕ_c 和 ϕ_h 都是相同的多层感知机(Multilayer Perceptron, MLP)。具体表现为使用Swish激活函数的双层MLP,第一层参数升至两倍,第二层参数下降为最初的大小。FGAT使用池化层提取节点特征的综合特征,并结合边缘特征得到更丰富的邻居节点的相关特征。最终,使用广义的张量积更新得到新的节点特征。FGAT与GAT的不同之处在于得到了节点与邻居节点之间更丰富的权重,同时避免多头注意力以及充分利用边缘特征。

3 实验与结果分析

3.1 实验设置

3.1.1 数据集及评价指标 本研究选用MoleculeNet的属性分类数据集^[10]以开展验证工作,且依据8:1:1的比例对数据集进行划分。为深入探究不同数据集划分方法对模型性能的影响,本研究引入两种划分方法:Scaffold Split和Random Split。其中,Scaffold Split^[20]是专为分子数据集设计的划分方法,其核心理念是将具有相同分子骨架的分子作为一个整体子集进行分割,这种划分方式更贴近真实世界中的预测任务场景,从而为模型性能的评估提供更具现实意义的参考依据^[21];Random Split是一种广泛应用于各类数据集划分的通用方法。本研究采用标准的Random Split方式对数据集进行处理,以一种更客观、中立的视角开展模型的消融实验,进而更全面地剖析模型在不同数据划分条件下的性能表现。关于所选用数据集的分子数量、预测任务以及详细描述见表2。

表2 数据集相关信息

Table 2 Dataset-related information

项目	数据集					
	BBBP	Tox21	ClinTox	HIV	BACE	SIDER
分子数目	2039	7831	1478	41127	1513	1427
任务数目	1	12	2	1	1	27
任务描述	标记药物分子是否具有血脑屏障渗透性	预测药物分子毒性,12个毒性靶点	预测FDA批准状态和毒性原因临床试验失败的药物	预测药物分子是否具有抑制HIV复制的活性	预测药物分子对AD关键酶(BACE-1)的抑制作用	预测药物分子的多重不良反应

在模型性能评价指标方面,本研究选择受试者操作特征曲线下面积(Area Under the ROC Curve, AUC)作为关键指标。根据相关文献和领域内的标准评估体系,当AUC值为0.5时,说明模型无实际区分能力,与随机猜测无异;当AUC值达到0.7及以上时,模型的区分能力相对可接受,通常可视为在特定应用情境下具有一定的预测效能。在具体实验操作中,针对单任务模型,设定10个固定随机种子,运行模型并取其平均值作为该任务模型的性能指标;对于多任务模型,分别对每个任务运行10个随机种子,先计算每个任务的平均值,再对所有任务的平均值进行二次平均,以此来综合评估多任务模型的整体性能。

3.1.2 训练细节 在模型训练环节,依据数据集的规模差异,灵活调整批次大小以优化训练效率。具体而言,针对HIV数据集,采用200作为批次大小;对于Tox21数据集,批次大小设定为32;而其他数据集则

统一采用8作为批次大小。鉴于多数数据集存在标签分布不均衡的情况,选用Focal Loss^[22]作为损失函数,以有效缓解类别不平衡对模型训练的影响,提升模型对少数类别的学习能力。除此之外,对所有数据集使用相同的训练策略。在训练过程中,模型共运行500轮次,未使用Dropout等正则化方法,使用Adam优化器(lr=1e-4)进行训练。

3.2 对比实验

为系统地评估Dual-FFTGCN模型的性能,将其与两类模型进行全面比较:一类是未经预训练的模型,包括GCN^[5]、GIN^[7]、GraphSAGE^[23]、SchNet^[24]和MGCN^[25];另一类是经过预训练的模型,如CHEM-BERT^[4]、MolCLR^[14]和MOCO^[15]。实验结果如表3所示,Dual-FFTGCN在所有数据集上均显著优于所有未经预训练的基准模型,并且与大多数预训练模型相比,依然展现出较强的竞争力。

表3 使用ROC-AUC作为评价指标的对比实验(%)
Table 3 Comparative experiments using ROC-AUC as the evaluation metric (%)

模型	数据集					
	BBBP	Tox21	ClinTox	HIV	BACE	SIDER
GCN	71.8±0.9	70.9±2.6	62.5±2.8	74.0±3.0	71.6±2.0	53.6±3.2
GraphSAGE	64.5±3.1	74.5±0.4	55.8±6.2	75.1±0.8	64.6±4.7	56.7±0.7
GIN	65.8±4.5	74.0±0.8	58.0±4.4	75.3±1.9	70.1±5.4	57.3±1.6
SchNet	84.8±2.2	77.2±2.3	71.5±3.7	70.2±3.4	77.6±1.1	53.9±3.7
MGCN	85.0±6.4	70.7±1.6	63.4±4.2	73.8±1.6	73.4±3.0	55.2±1.8
CHEM-BERT	72.4±0.9	77.4±0.5	99.0±0.3	77.6±1.6	82.0±1.7	63.1±0.6
MolCLR	73.6±0.5	79.8±0.7	93.2±1.7	80.6±1.1	89.0±0.3	68.0±1.1
MOCO	71.6±1.0	76.7±0.4	81.6±3.7	78.3±0.4	82.6±0.3	61.2±0.6
Dual-FFTGCN	85.6±2.6	79.6±0.5	75.0±1.7	81.2±0.6	84.1±1.2	68.4±1.2

Dual-FFTGCN通过巧妙地融合频域和空间域的信息,显著提升分子属性预测的性能。在多个标准数据集上的实验结果表明,该模型不仅在准确率上优于或与现有的基准模型相当,而且在处理复杂的分子结构时,展现出更强的泛化能力和预测性能,这为药物发现领域提供一种高效且可靠的解决方案,有望在实际应用中发挥重要作用。

3.3 消融实验

为深入理解Dual-FFTGCN模型中各个组件的贡献,设计一系列消融实验,系统地移除模型中的不同部分,并评估其对整体性能的影响。具体来说,考虑4种模型配置:移除图注意力网络(del_FGAT)、移除图卷积网络(del_GCN)、移除池化层(del_Pool)以及完整模型(FULL)。实验结果如图2所示,涵盖6个标准数据集:

HIV、BACE、BBBP、ClinTox、Tox21和SIDER。

实验结果表明,Dual-FFTGCN模型中的各个组件,包括图注意力网络、图卷积网络和池化层,均对模型的性能产生显著影响。特别是,移除池化层后模型性能接近表现最好的单卷积核图卷积神经网络,这表明了融合两种卷积核图神经网络方法的有效性。同时,模型在添加池化层后,无关数据集大小均能提高模型的相关性能,这说明了基于傅里叶滤波器的图池化层具备效能的同时具有普适性,这些发现为进一步优化和改进Dual-FFTGCN模型提供重要参考。

4 讨论

本研究探索不同卷积核图神经网络的融合方

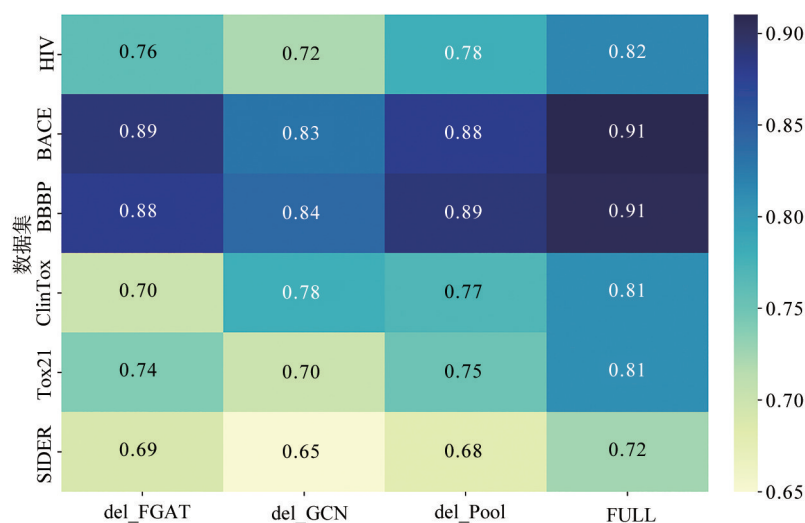


图2 各数据集上的消融实验热力图

Figure 2 Heatmap of ablation experiments on various datasets

法,通过将无权图送入基于谱域的图神经网络及将有权图送入基于空域的图神经网络,提高对于分子图的利用。同时,使用基于傅里叶变换的滤波器成功提高图神经网络对于分子结构的表征能力。更值得注意的是本文模型未经过大规模的预训练任务,仅通过对数据集的训练和有限的先验知识就达到了一系列经典的预训练模型的性能。

在实验过程中发现模型在小型数据集上的表现未能达到预期水平,这一现象在一定程度上凸显了预训练任务实施的必要性与迫切性。小型数据集往往难以提供足够的信息用于模型的有效训练,导致模型在学习数据特征和模式时受到限制,从而影响其性能表现。未来将有计划对自监督任务、对比学习等预训练任务和数据增强技术进行研究,使得模型更稳定且性能更进一步,从而促进人工智能辅助药物开发和药物发现。

【参考文献】

- [1] Deng JA, Yang ZB, Ojima I, et al. Artificial intelligence in drug discovery: applications and techniques[J]. *Brief Bioinform*, 2022, 23(1): bbab430.
- [2] Zhang Y, Luo MQ, Wu P, et al. Application of computational biology and artificial intelligence in drug design[J]. *Int J Mol Sci*, 2022, 23(21): 13568.
- [3] Yang NZ, Wu HJ, Zeng KP, et al. Molecule generation for drug design: a graph learning perspective[J/OL]. *Fundam Res*. (2024-12-20). <https://doi.org/10.1016/j.fmr.2024.11.027>.
- [4] Kim H, Lee J, Ahn S, et al. A merged molecular representation learning for molecular properties prediction with a web-based service[J]. *Sci Rep*, 2021, 11(1): 11028.
- [5] Kipf TN, Welling M. Semi-supervised classification with graph convolutional networks[EB/OL]. (2017-02-22). <https://arxiv.org/abs/1609.02907>.
- [6] Veličković P, Cucurull G, Casanova A, et al. Graph attention networks[EB/OL]. (2018-02-04). <https://arxiv.org/abs/1710.10903>.
- [7] Xu KY, Hu WH, Leskovec J, et al. How powerful are graph neural networks?[EB/OL]. (2019-02-22). <https://arxiv.org/abs/1810.00826>.
- [8] Xia J, Zhang LC, Zhu X, et al. Understanding the limitations of deep models for molecular property prediction: insights and solutions[C]//

Proceedings of the 37th International Conference on Neural Information Processing Systems. Red Hook, NY, USA: Curran Associates Inc., 2023: 64774-64792.

- [9] Yi K, Zhang Q, Fan W, et al. FourierGNN: rethinking multivariate time series forecasting from a pure graph perspective [C]// Proceedings of the 37th International Conference on Neural Information Processing Systems. Red Hook, NY, USA: Curran Associates Inc., 2023: 69638-69660.
- [10] Wu ZQ, Ramsundar B, Feinberg EN, et al. MoleculeNet: a benchmark for molecular machine learning[J]. *Chem Sci*, 2018, 9(2): 513-530.
- [11] Devlin J, Chang MW, Lee K, et al. BERT: pre-training of deep bidirectional transformers for language understanding [C]// Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Stroudsburg, PA, USA: ACL, 2019: 4171-4186.
- [12] Batatia I, Kovács DP, Simm GN, et al. MACE: higher order equivariant message passing neural networks for fast and accurate force fields [C]// Proceedings of the 36th International Conference on Neural Information Processing Systems. Red Hook, NY, USA: Curran Associates Inc., 2022: 11423-11436.
- [13] Du WT, Du YQ, Wang LM, et al. A new perspective on building efficient and expressive 3D equivariant graph neural networks [C]// Proceedings of the 37th International Conference on Neural Information Processing Systems. Red Hook, NY, USA: Curran Associates Inc., 2023: 66647-66674.
- [14] Wang YY, Wang JR, Cao ZL, et al. Molecular contrastive learning of representations via graph neural networks[J]. *Nat Mach Intell*, 2022, 4(3): 279-287.
- [15] Zhu YQ, Chen DS, Du YQ, et al. Molecular contrastive pretraining with collaborative featurizations[J]. *J Chem Inf Model*, 2024, 64(4): 1112-1122.
- [16] Yang YX, Wang ZX, Ahadian P, et al. A deep multimodal representation learning framework for accurate molecular properties prediction [C]// Proceedings of the Great Lakes Symposium on VLSI 2024. New York, NY, USA: Association for Computing Machinery, 2024: 760-765.
- [17] Quesado P, Torres LH, Ribeiro B, et al. A hybrid GNN approach for improved molecular property prediction[J]. *J Comput Biol*, 2024, 31(11): 1146-1157.
- [18] Chung J, Gulcehre C, Cho KH, et al. Empirical evaluation of gated recurrent neural networks on sequence modeling[EB/OL]. (2014-12-11). <https://arxiv.org/abs/1412.3555>.
- [19] Satorras VG, Hoogeboom E, Welling M. E(n) equivariant graph neural networks[C]// Proceedings of the 38th International Conference on Machine Learning. Chia Laguna Resort, Sardinia, Italy: PMLR, 2021: 9323-9332.
- [20] Ramsundar B, Liu BW, Wu ZQ, et al. Is multitask deep learning practical for pharma?[J]. *J Chem Inf Model*, 2017, 57(8): 2068-2076.
- [21] Yang K, Swanson K, Jin W, et al. Analyzing learned molecular representations for property prediction[J]. *J Chem Inf Model*, 2019, 59(8): 3370-3388.
- [22] Lin TY, Goyal P, Girshick R, et al. Focal loss for dense object detection [C]// 2017 IEEE International Conference on Computer Vision (ICCV). Piscataway, NJ, USA: IEEE, 2017: 2999-3007.
- [23] Hamilton WL, Ying R, Leskovec J. Inductive representation learning on large graphs[C]// Proceedings of the 31st International Conference on Neural Information Processing Systems. Red Hook, NY, USA: Curran Associates Inc., 2017: 1025-1035.
- [24] Schütt KT, Sauceda HE, Kindermans PJ, et al. SchNet-a deep learning architecture for molecules and materials[J]. *J Chem Phys*, 2018, 148(24): 241722.
- [25] Lu CQ, Liu Q, Wang C, et al. Molecular property prediction: a multilevel quantum interactions modeling perspective [C]// Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence. Palo Alto, CA, USA: AAAI Press, 2019: 1052-1060.

(编辑:谭斯允)