

基于双路TCN的长短距离融合学习转录因子结合位点预测

吴志强¹, 宋佳智¹, 姜静清¹, 罗蕊²

(1. 内蒙古民族大学 计算机科学与技术学院; 2. 内蒙古民族大学 生命科学与食品学院, 内蒙古 通辽 028000)

摘要: 准确预测DNA与转录因子的结合位点对深入理解基因表达及调控机理具有重要意义。卷积神经网络(CNN)和长短时记忆网络(LSTM)已成功应用于DNA-转录因子结合位点预测任务, 准确性相比传统机器学习方法提升明显。然而, CNN仅擅长学习局部空间信息, 无法建模DNA序列中的长距离依赖关系; LSTM网络的顺序处理特性无法实现并行运算, 计算效率偏低。为此, 提出一种结合双路时间卷积网络(TCN)和长短距离融合学习机制的模型解决上述问题。在网络结构层面, 使用时间卷积网络作为序列特征提取器, 兼具长距离建模和并行处理的优势, 而且双路结构使模型能分离学习DNA互补特征, 一定程度上提升了模型的稳定性; 在特征学习层面, 充分利用时间卷积网络不同层的上下文信息建模能力设计了长短距离融合学习策略, 增强了预测特征的代表能力。在165个ChIP-seq数据集上的实验结果表明, 该方法的各项指标结果优于当前流行的基于深度学习的方法。通过利用具有不同距离依赖信息的时序特征, 可为转录因子结合位点预测提供一个有价值的框架。

关键词: 转录因子; 结合位点; 序列数据处理; 时间卷积网络; 特征融合

DOI: 10.11907/rjdk.241671

开放科学(资源服务)标识码(OSID):



中图分类号: TP18

文献标识码: A

文章编号: 1672-7800(2025)003-0031-06

Transcription Factor Binding Sites Prediction with Long-Short Distance Fusion Learning Based on Dual-path TCN

WU Zhiqiang¹, SONG Jiazhi¹, JIANG Jingqing¹, LUO Rui²

(1. College of Computer Science and Technology, Inner Mongolia Minzu University;

2. College of Life Sciences and Food Engineering, Inner Mongolia Minzu University, Tongliao 028000, China)

Abstract: Accurate identification of transcription factor binding sites (TFBSs) is crucial for understanding gene expression and regulatory mechanisms. Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM) models have significantly improved accuracy in this task compared to traditional machine learning approaches. However, CNNs specialize in learning local spatial features but ignore long-distance dependencies in DNA sequences, while LSTMs are proficient in learning sequential relationships but computationally inefficient due to a lack of parallel computing ability. This paper proposes a novel dual-path sequential network integrating long-short distance fusion learning to address the above issues. In terms of structure, this paper employs the Temporal Convolutional Network (TCN) as the feature extractor that supports sequential modeling and parallel processing. The dual-path structure can learn complementary DNA features, improving the learning stability. In terms of features, this paper leverages the context information modeling capability of TCN, and designs a long-short distance fusion learning strategy to strengthen the feature representation for prediction. The experiment results on 165 ChIP-seq datasets show that our

收稿日期: 2024-08-07

扫描二维码阅读全文:



基金项目: 国家自然科学基金项目(62162050); 内蒙古民族大学博士科研启动基金项目(KYQD23006, BS672); 内蒙古自然科学基金项目(2021BS03036); 蓖麻产业技术创新内蒙古自治区工程研究中心开放课题(MDK2021004, MDK2023012); 内蒙古自治区蓖麻产业协同创新中心开放课题(MDK2022016)

作者简介: 吴志强(1986-), 男, 内蒙古民族大学计算机科学与技术学院硕士研究生, 研究方向为生物医学信息处理; 宋佳智(1993-), 男, 博士, 内蒙古民族大学计算机科学与技术学院讲师, 研究方向为人工智能与生物信息学; 姜静清(1968-), 女, 博士, 内蒙古民族大学计算机科学与技术学院教授, 研究方向为生物信息学; 罗蕊(1988-), 女, 博士, 内蒙古民族大学生命科学与食品学院讲师, 研究方向为细胞生物学。本文通讯作者: 姜静清。

method outperforms the popular deep learning based methods. This study introduces a valuable framework for TFBSs prediction by combining sequential features with different distance dependency information.

Key Words: transcription factor; binding site; sequence data processing; temporal convolutional network; feature fusion

0 引言

转录因子通过结合特定的DNA序列片段调控下游基因的表达,这些特定的DNA片段称为结合位点^[1]。确定转录因子的结合位点能够帮助研究人员更好地理解基因内部的调控机制及细胞功能,因此在基因治疗和药物设计等生物医药领域具有重要的应用价值。

近年来,高通量测序技术的不断发展产生了海量的转录因子结合位点数据。如何通过已有的计算方法挖掘数据中的隐含模式,将其应用至新序列的结合位点预测成为关键问题。早期研究人员通常使用支持向量机^[2]、隐马尔可夫模型^[3]等传统机器学习模型建模DNA序列与结合位点的关系。然而,由于模型自身的表达能力受限,传统机器学习方法面对日益增加的DNA序列数据难以进一步提升准确率。

随着深度学习方法的兴起,计算机视觉和自然语言处理技术得到了迅速发展。在转录因子结合位点预测领域,研究人员也逐步开始应用深度学习模型提取序列特征以提升预测效果^[4]。常用的深度学习模型主要包括两类:卷积神经网络(Convolutional Neural Network, CNN)和循环神经网络。如DeepBind^[5]、DeepSEA^[6]和CNN-zeng^[7]等方法应用CNN构建DNA序列中的局部空间关系,以获取更丰富的序列表征模式;DeepSite^[8]和DanQ^[9]则通过长短时记忆网络(Long Short-Term Memory, LSTM)学习DNA双链中的序列关系,增强特征的鲁棒性。得益于模型强大的特征表示能力,上述方法在各项评价指标上的表现远优于传统机器学习方法。

尽管上述方法的预测性能远优于传统方法,但在DNA序列特征提取方面尚存在以下缺陷:①CNN能获取特征中的局部空间模式,但无法学习序列远端特征间的依赖关系;②循环神经网络适合学习DNA数据中的隐含序列关系,但自身具备的顺序处理特性使其在处理大规模数据时面临内存占用过高和训练效率低下的问题。

为解决上述问题,本文提出一种全新的双路时间卷积网络(Temporal Convolutional Network, TCN)框架进行转录

因子结合位点预测。该框架具有三方面优势:①采用的时间卷积网络不仅能学习序列中不同距离特征间的全局依赖关系,而且具有并行计算优势,计算效率高;②框架设计了双路网络结构分离提取DNA双链的序列特征,可提升特征学习的稳定性;③提出一种长短距离融合学习策略,能够有效联合时间卷积网络不同层的特征,增强网络对DNA序列的表征能力。在165个ENCODE ChIP-seq数据集上的大量实验结果表明,本文方法与基于卷积神经网络和循环神经网络等方法相比,在分类精度等评测指标上具有更好的表现,充分验证了本文方法的有效性。

1 材料与方

1.1 数据集

本文使用DNA百科全书项目(Encyclopedia of DNA Elements, ENCODE)^[10]产生的ChIP-seq数据集,数据集的设定策略与D-SSCA^[11]和DSAC^[12]完全一致:165个子数据集参与实验,包含29种不同类型的转录因子结合位点数据;每一个子数据集的训练样本和测试样本占比分别为80%与20%,训练样本进一步随机划分出20%作为验证集以优化模型超参数。

1.2 模型架构

双路时间卷积网络框架如图1所示。该框架以互补的DNA序列作为输入,One-hot编码层将序列中的碱基对表征为数字向量,并输入双路时间卷积网络提取特征。所提取的特征通过长短距离融合学习策略形成更具表达力特征,该特征经过最后的分类层输出DNA序列为转录因子结合位点的概率。

(1)互补的DNA序列。生物体中的DNA序列由两条单链组成,其根据固定的碱基互补原则配对(腺嘌呤A与胸腺嘧啶T配对,鸟嘌呤G与胞嘧啶C配对)形成双链结构。ChIP-seq数据集中每个样本的DNA序列仅包含101个碱基。本文根据互补原则为每一个样本生成标签一致的配对序列,通过扩充训练样本数量增加互补对序列信息,以提升特征学习的稳定性及模型预测的准确性。

(2)One-hot编码层。令DNA序列表示为S =

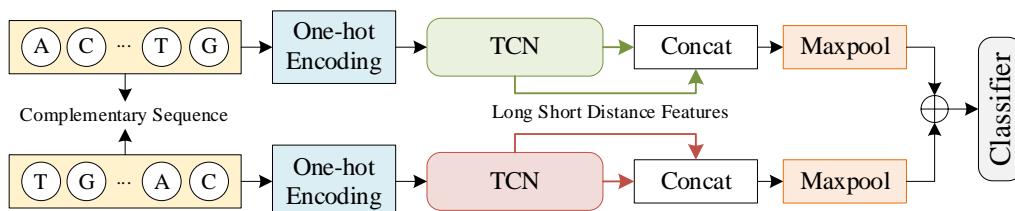


Fig. 1 The framework of dual-path TCN

图1 双路时间卷积网络框架

(s_1, s_2, \dots, s_n) , 其中 n 为序列长度, $s_i (1 \leq i \leq n)$ 属于碱基集合 $\{A, C, G, T\}$ 。one-hot 编码层将 A、C、G、T 分别表示为向量 $[1, 0, 0, 0]^T$ 、 $[0, 1, 0, 0]^T$ 、 $[0, 0, 1, 0]^T$ 、 $[0, 0, 0, 1]^T$ 。因此, DNA 序列 S 经过该编码层最终可表示为大小为 $4 \times n$ 的矩阵。

(3) 时间卷积网络。为提取序列数据中的时序特征, 研究人员通常使用 LSTM 作为特征提取器^[8-9]。但是, 该网络只能按照顺序处理序列数据, 序列越长计算效率越低。为解决该问题, 研究人员提出时间卷积网络(TCN)进行序列建模任务^[13]。因果卷积模块是 TCN 的重要组成部分,

其遵守严格的时间约束机制, 还兼具序列的并行处理能力。由于 DNA 序列中隐含更复杂的先后依赖关系, 本文采用一种 TCN 变体提取 DNA 序列特征^[14]。该结构摒弃了原始 TCN 中的因果约束, 相同条件下输出的特征包含更广泛的序列关系。同时, 本文融合了各级中间特征以提升模型的预测性能。

时间卷积网络由 1 个 1×1 卷积层和 6 个膨胀残差卷积层组成。 1×1 卷积调整输入数据与膨胀残差卷积输入特征之间的通道关系, 多层膨胀残差卷积逐层建模表示序列特征。时间卷积网络架构如图 2 所示。

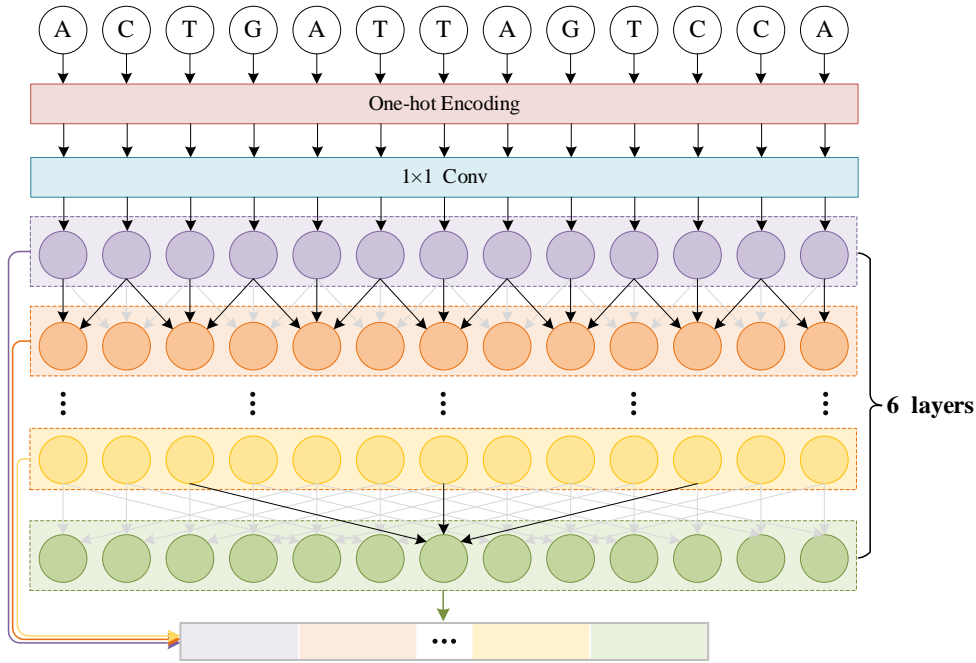


Fig. 2 The structure of TCN

图 2 时间卷积网络架构

膨胀残差卷积模块是时间卷积网络的主要组成部分, 其包含膨胀卷积、ReLU 激活函数、 1×1 卷积、Dropout 层和残差连接。具体结构如图 3 所示。

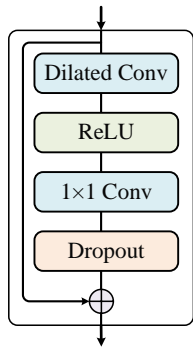


Fig. 3 The structure of dilated residual convolutional module

图 3 膨胀残差卷积模块结构

相较于普通卷积, 膨胀卷积具有更大的感受野, 能在序列中建模远距离特征间的关系。其计算公式如下^[13]:

$$G(s) = \sum_{i=-(k-1)/2}^{(k-1)/2} g_i \cdot x_{s+d \cdot i}, \quad (1)$$

其中, $G(\cdot)$ 表示作用在序列中元素 s 上的膨胀卷积操作, g 表示卷积核, k 表示卷积核大小 (奇数), x 表示输入的 DNA 序列, $s + d \cdot i$ 隐含了计算方向。 d 表示膨胀系数, 本文 d 以 2^i (i 从 0 开始) 的规律逐层递增。添加 ReLU 激活函数和 Dropout 层是为了增加网络的表达能力并防止过拟合现象。同时, 为防止网络过深导致的梯度消失等问题, 本文使用残差结构衔接每一层的输入和输出, 因而要求输入序列与输出序列的长度和通道数相等。为此, 本文将膨胀卷积核大小设置为 3, 序列的填充数与膨胀系数设置相同。 1×1 卷积层的输入输出通道数保持一致, 其步长为 1。

(1) 长短距离特征融合学习策略。由于 TCN 最后一层的输出特征中只能建模输入序列的远端关系, 无法涵盖近距离特征之间的关系表达。因此, 本文提出一种结合长短距离特征融合学习策略增强模型的特征表达力。该策略执行过程分为两部分: 长短距离特征合并 (图 1 中的 “concat”) 和双链特征融合 (图 1 中的 “ \oplus ”)。

采用长短距离特征合并策略基于以下事实: TCN 不同层特征表示不同距离的依赖关系, 网络越深, 建模的特征

依赖关系距离越长,如图2所示。因此,为充分利用不同层的特征,本文合并了各层特征参与最终的预测过程。双链特征融合策略的目的是控制计算复杂度,同时保留特征中的关键信息。为此,本文将合并后的特征通过最大池化层筛选突出特征,并将双链特征相加作为最终的预测特征。该策略保留了原始序列的关键信息,与常用的串联特征方式相比也能一定程度地减少计算成本。实验结果验证了长短距离特征融合学习策略的有效性。

(2)分类器。分类器由一层全连接网络和Sigmoid激活函数组成。经过Sigmoid函数后,使用二元交叉熵函数计算最终损失:

$$L = \frac{1}{N} \sum_{j=1}^N \bar{y}_j \times \log(y_j) + (1 - \bar{y}_j) \times \log(1 - y_j) \quad (2)$$

其中, L 表示 N 个样本计算得到的损失函数值。 \bar{y}_j 表示第 j 个样本的标签,取值为0或1。当 $\bar{y}_j = 1$ 时,表示当前样本 j 是结合位点,反之则不是。 y_j 表示第 j 个样本预测为结合位点的概率($0 \leq y_j \leq 1$)。

2 实验结果与分析

2.1 超参数设定

本文使用PyTorch深度学习框架实现模型。为获得更优的效果,通过比较模型在验证集上的性能调整超参数,表1展示了主要的超参数及其对应的搜索空间和最优值。具体如下:使用Adam优化器更新模型,初始学习率设置为0.001^[15]。为更精细地优化模型,通过学习率衰减策略调整学习率大小,即每训练10轮,学习率减少为原来的1/5。每个子数据集均训练40轮,设置批量大小为128。在双路TCN中,膨胀残差卷积层的输入通道和输出通道数均为64。

Table 1 Hyperparameters of the proposed method and the corresponding search space and optimal value

表1 本文方法的超参数与对应搜索空间及最优值

超参数	搜索范围	最优值
学习率	[0.000 5, 0.003] ^(a:搜索步长为5e-4)	0.001
批量大小	{32, 64, 128, 256}	128
优化器	{Adam, SGD}	Adam
膨胀残差卷积输入通道与输出通道大小	{32, 64, 128}	64

2.2 评价指标

本文采用的评价指标与文献[11]、[12]一致,使用准确率(Accuracy, ACC)、接受者操作特征(Receiver Operating Characteristic, ROC)、曲线下的面积(Area under the ROC Curve, ROC-AUC)和精确率—召回率(Precision-Recall, PR)曲线下的面积(Area under the PR Curve, PR-AUC)测定和比较不同方法的性能。

在二分类任务中,为了衡量模型性能,通常将预测结

果分为4类:真阳性(True Positive, TP),表示真实值和预测值都为正样本;真阴性(True Negative, TN),表示真实值和预测值都为负样本;假阳性(False Positive, FP),即真实值为负样本,但预测值为正样本;假阴性(False Negative, FN),即真实值为正样本,但预测值为负样本。本文采用的模型评价指标与4类指标紧密关联。其中,ACC计算分类正确的样本数占总样本数的比例。计算公式如下:

$$ACC = \frac{TP + TN}{TP + TN + FP + FN} \quad (3)$$

考虑到165个ChIP-seq数据集存在正负样本不均衡的情况,因此进一步通过ROC-AUC和PR-AUC测量模型性能^[11-12]。ROC曲线通过绘制真阳性率(True Positive Rate, TPR)与假阳性率(False Positive Rate, FPR)评估模型性能。其中,TPR计算模型预测为正确的正样本在所有正样本中的比例。计算公式如下:

$$TPR = \frac{TP}{TP + FN} \quad (4)$$

FPR计算模型预测为正样本,但实际为负样本在所有负样本中的比例。计算公式如下:

$$FPR = \frac{FP}{FP + TN} \quad (5)$$

精确率(Precision)计算模型预测为正确的正样本在所有预测的正样本中的比例。计算公式如下:

$$Precision = \frac{TP}{TP + FP} \quad (6)$$

召回率(Recall)与TPR的含义相同,计算过程如式(4)所示。

2.3 对比结果

为验证本文方法的性能,将其与目前流行的基于深度学习的方法进行对比,对比方法包括DeepBind(CNN)^[5]、DanQ(CNN+LSTM)^[9]、CRPTS(CNN+LSTM)^[16]、DLBSS(CNN)^[17]、D-SSCA(CNN)^[11]、DSAC(CNN)^[12]和DeepSTF(Transformer)^[18]。值得一提的是,DeepSTF同时使用了DNA序列信息和形状信息进行转录因子结合位点预测。所有方法均在165个ChIP-seq数据集上进行实验,结果如图4、表2所示。在图4中,长方形框内的中线表示中位数,上下两端表示上四分位数和下四分位数;长方形外的两端表示各个指标的最小值和最大值;菱形标记代表异常值。

如图4所示,除最大值与DeepSTF等方法相近外,本文方法在3项指标的最小值和中值等均高于其它方法。从全局上看,本文方法各项指标结果分布在更高的值区间。由表2可知,本文方法在测试数据集上的平均ACC、ROC-AUC与PR-AUC分别为0.826、0.902和0.906,与次优方法DSAC相比分别提升了1.0%、1.5%和1.5%。这一定程度上表明本文方法相较于对比方法具备更强的泛化能力,主要原因如下:①本文方法使用了擅长序列信息处理的TCN作为基础特征提取器,与CNN和LSTM相比,TCN能在序列并行处理的同时有效建模DNA远端特征间的关系;②本文设计了双路TCN网络分别提取双链DNA特征,通过组

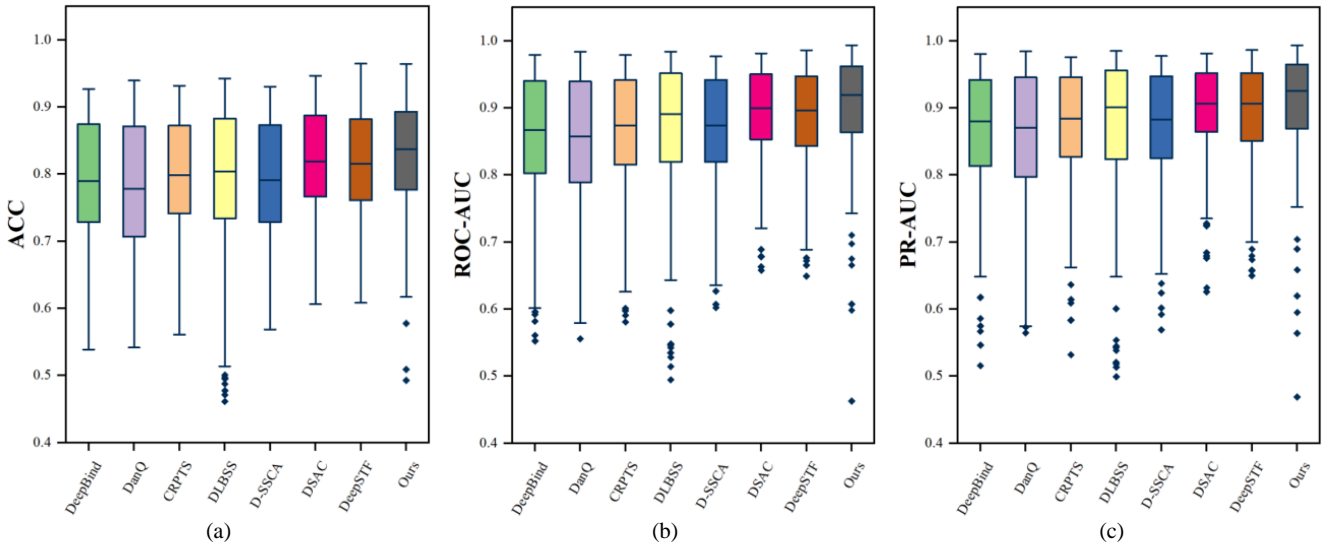


Fig. 4 Comparison of ACC, ROC-AUC and PR-AUC of different methods on 165 ChIP-seq datasets

图 4 不同方法在 165 个 ChIP-seq 数据集上的 ACC、ROC-AUC 与 PR-AUC 比较

Table 2 Comparison of the average ACC, ROC-AUC and PR-AUC of different methods on 165 ChIP-seq datasets

表 2 不同方法在 165 个 ChIP-seq 数据集上的 ACC、ROC-AUC 与 PR-AUC 均值比较

Method	ACC	ROC-AUC	PR-AUC
DeepBind ^[5]	0.785	0.853	0.858
DanQ ^[9]	0.782	0.849	0.855
CRP-TS ^[16]	0.793	0.862	0.867
DLBSS ^[17]	0.793	0.865	0.871
D-SSCA ^[11]	0.793	0.867	0.871
DSAC ^[12]	0.816	0.887	0.891
DeepSTF ^[18]	0.814	0.883	0.890
本文方法	0.826	0.902	0.906

合不同性质的特征提升模型的表达能力;③本文提出的长短距离特征融合策略通过联合网络各级特征增强结果的特征表示能力,提升了模型性能。值得一提的是,对比方法 DeepSTF 虽然使用了 DNA 序列和形状两种异质数据,但由于 Transformer 网络的引入,在本文训练数据有限的场景中难以发挥其强大的复杂关系学习能力,导致 DeepSTF 的性能难以达到最优。

2.4 消融实验

在模型构建过程中,为选择模型结构和预测特征,从输入数据、网络结构和预测特征 3 个视角分别进行实验分析。在输入数据层面,使用原始的 DNA 单链序列、互补序列及双链序列输入单路 TCN 网络验证不同输入数据的效果。在网络结构层面,本文以 DNA 双链序列作为输入数据,着重对比单路 TCN 及结构一致的双路 TCN 模型的性能。上述两组实验仅使用 TCN 输出的最后一层特征进行分类。在预测特征层面,本文将 DNA 双链序列输入至双路 TCN 网络,分别评估使用最后一层特征和合并的长短距离特征作为最终预测特征的效果。所有实验均在 165 个 ChIP-seq 数据集上进行,超参数设定与 2.1 节的描述一致。

2.4.1 输入数据

直观上,使用单链特征与双链特征的不同之处在于后者的训练数据量是前者的两倍,单链特征蕴含的信息量一定程度上比双链特征更少。因此,使用双链特征时的预测准确率将优于仅使用单链特征。实验结果也证明了这一猜想的有效性。如表 3 所示,仅使用原始特征和互补特征的 ACC、ROC-AUC 及 PR-AUC 几乎相等,差值仅在 0.2% 以内。当使用双链特征时,3 项指标比仅使用单链特征均有所提升,表明使用信息含量更多的双链特征能一定程度上提升模型的泛化能力。

Table 3 Performance comparison of using single-chain feature and double-chain feature

表 3 使用单链特征与双链特征结果比较

Method	ACC	ROC-AUC	PR-AUC
仅原始特征	0.810	0.887	0.892
仅互补特征	0.811	0.889	0.894
双链特征	0.816	0.891	0.896

2.4.2 网络结构

相较于单路 TCN 模型,双路 TCN 通过分离提取 DNA 双链特征,能一定程度上增加特征学习的稳定性。如表 4 所示,实验结果表明使用双路 TCN 模型的预测效果更佳,ACC、ROC-AUC 及 PR-AUC 相比使用单路网络分别提高了 0.3%、0.7% 和 0.6%。说明在本文的任务中,双路结构相较于单路结构更优。

2.4.3 预测特征

与仅使用最后一层特征作为最终预测特征相比,合并

Table 4 Performance comparison of using single-path network and dual-path network

表 4 使用单路网络与双路网络结果比较

Method	ACC	ROC-AUC	PR-AUC
单路网络	0.816	0.891	0.896
双路网络	0.819	0.898	0.902

TCN各层特征能够有效地利用不同距离特征间的依赖关系,从而提升预测特征的代表能力。如表5所示,模型合并特征后的ACC、ROC-AUC与PR-AUC分别达到了82.6%、90.2%和90.6%,3项指标均比仅使用最后一层特征表现优异,充分验证了采用长短距离特征融合策略的有效性。

Table 5 Performance comparison of using last-layer feature of TCN and the long-short distance fusion feature

表5 使用TCN最后一层特征与长短距离融合特征结果比较

Method	ACC	ROC-AUC	PR-AUC
最后一层特征	0.819	0.898	0.902
合并的长短距离特征	0.826	0.902	0.906

3 结语

本文提出一种双路时间卷积网络模型和长短距离特征融合学习策略完成DNA-转录因子结合位点预测。该模型采用时间卷积网络作为DNA序列的特征提取器,与常用的CNN和LSTM网络相比,时间卷积网络具有长距离建模和并行运算等优势。同时,设计的双路网络结构使模型能分离提取DNA双链特征,从而提升了特征学习的稳定性。此外,本文提出了长短距离特征融合学习策略,该策略充分利用了时间卷积网络的上下文信息建模能力,显著增强了最终预测特征的代表能力。在165个ChIP-seq数据集上的实验结果表明,本文方法相较于现有的基于CNN、基于LSTM和混合模型等深度学习方法,在ACC、ROC-AUC及PR-AUC 3项指标上均取得了更优的结果,充分表明本文方法在DNA-转录因子结合位点预测任务中的有效性。

在实验结果中,本文设计的双路模型性能虽然比单路模型更优,但成倍的网络参数带来的性能提升并不明显。因此,在未来的工作中,本文将继续研究并设计更高效的深度学习特征提取模型,以进一步提升DNA-转录因子结合位点预测的准确性。

参考文献:

- [1] JIANG B W, FENG Z J, HUANG W H. DNA transcription factor binding site prediction based on split-attention mechanism[J]. *Software Guide*, 2024, 23(2): 32-39.
姜博文, 冯子健, 黄伟鸿. 基于分裂注意力机制的DNA转录因子结合位点预测[J]. *软件导刊*, 2024, 23(2): 32-39.
- [2] HU J L, WANG J R, LIN J A, et al. MD-SVM: a novel SVM-based algorithm for the motif discovery of transcription factor binding sites[J]. *BMC Bioinformatics*, 2019, 20: 200.
- [3] WON K J, REN B, WANG W. Genome-wide prediction of transcription factor binding sites using an integrated model[J]. *Genome Biology*, 2010, 11(1): 1-17.
- [4] HE Y, SHEN Z, ZHANG Q, et al. A survey on deep learning in DNA/RNA motif mining[J]. *Briefings in Bioinformatics*, 2021, 22(4): 1-10.
- [5] ALIPANAHI B, DELONG A, WEIRAUCH M T, et al. Predicting the sequence specificities of DNA-and RNA-binding proteins by deep learning[J]. *Nature Biotechnology*, 2015, 33(8): 831-838.
- [6] ZHOU J, TROYANSKAYA O G. Predicting effects of noncoding variants with deep learning-based sequence model[J]. *Nature Methods*, 2015, 12(10): 931-934.
- [7] ZENG H, EDWARDS M D, LIU G, et al. Convolutional neural network architectures for predicting DNA-protein binding[J]. *Bioinformatics*, 2016, 32(12): i121-i127.
- [8] ZHANG Y, QIAO S, JI S, et al. DeepSite: bidirectional LSTM and CNN models for predicting DNA-protein binding[J]. *International Journal of Machine Learning and Cybernetics*, 2020, 11: 841-851.
- [9] QUANG D, XIE X. DanQ: a hybrid convolutional and recurrent deep neural network for quantifying the function of DNA sequences[J]. *Nucleic Acids Research*, 2016, 44(11): e107.
- [10] ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome[J]. *Nature*, 2012, 489: 57-74.
- [11] ZHANG Y, WANG Z, ZENG Y, et al. A novel convolution attention model for predicting transcription factor binding sites by combination of sequence and shape[J]. *Briefings in Bioinformatics*, 2022, 23(1): 1-12.
- [12] YU Y, DING P, GAO H, et al. Cooperation of local features and global representations by a dual-branch network for transcription factor binding sites prediction[J]. *Briefings in Bioinformatics*, 2023, 24(2): 1-9.
- [13] BAI S, KOLTER J Z, KOLTUN V. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling[DB/OL]. <https://arxiv.org/abs/1803.01271>. pdf.
- [14] FARHA Y A, GALL J. MS-TCN: multi-stage temporal convolutional network for action segmentation[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019: 3575-3584.
- [15] KINGMA D P, BA J. Adam: a method for stochastic optimization[DB/OL]. <https://arxiv.org/abs/1412.6980>. pdf.
- [16] WANG S, ZHANG Q, SHEN Z, et al. Predicting transcription factor binding sites using DNA shape features based on shared hybrid deep learning architecture[J]. *Molecular Therapy-Nucleic Acids*, 2021, 24: 154-163.
- [17] ZHANG Q, SHEN Z, HUANG D S. Predicting in-vitro transcription factor binding sites using DNA sequence+ shape[J]. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 2019, 18(2): 667-676.
- [18] DING P, WANG Y, ZHANG X, et al. DeepSTF: predicting transcription factor binding sites by interpretable deep neural networks combining sequence and shape[J]. *Briefings in Bioinformatics*, 2023, 24(4): 1-11.

(责任编辑:黄健)