

冯国富, 卢胜涛, 陈 明, 等. 基于自注意力机制和改进的 K -BiLSTM 的水产养殖水体溶解氧含量预测模型[J]. 江苏农业学报, 2024, 40(3): 490-499.

doi: 10.3969/j.issn.1000-4440.2024.03.011

基于自注意力机制和改进的 K -BiLSTM 的水产养殖水体溶解氧含量预测模型

冯国富^{1,2}, 卢胜涛^{1,2}, 陈 明^{1,2}, 王耀辉³

(1. 上海海洋大学信息学院, 上海 201306; 2. 农业农村部渔业信息重点实验室, 上海 201306; 3. 南通龙洋水产有限公司, 江苏南通 226634)

摘要: 为精确预测水产养殖水体溶解氧含量, 本研究提出一种基于自注意力机制(ATTN)和改进的 K -means 聚类-基于残差和批标准化(BN)的双向长短期记忆网络(BiLSTM)的水产养殖水体溶解氧含量预测模型。首先, 根据环境数据的相似性, 使用改进的 K -means 算法将数据划分成若干类别; 然后, 在 BiLSTM 基础上构建残差连接和加入 BN 完成高层次特征提取, 利用 BiLSTM 的长期记忆能力保存特征信息; 最后, 引入自注意力机制突出不同时间节点数据特征的重要性, 进一步提升模型的性能。试验结果表明, 本研究提出的基于自注意力机制和改进的 K -BiLSTM 模型的平均绝对误差为 0.238、均方根误差为 0.322、平均绝对百分比误差为 0.035, 与单一的 BP 模型、CNN-LSTM 模型、传统的 K -means-基于残差和 BN 的 BiLSTM-ATTN 等模型相比具有更优的预测性能和泛化能力。

关键词: 水产养殖; 溶解氧预测; K -means 聚类; 双向长短期记忆网络(BiLSTM); 自注意力机制

中图分类号: TP391 **文献标识码:** A **文章编号:** 1000-4440(2024)03-0490-10

Prediction model of dissolved oxygen content in aquaculture water based on self-attention mechanism and improved K -BiLSTM

FENG Guo-fu^{1,2}, LU Sheng-tao^{1,2}, CHEN Ming^{1,2}, WANG Yao-hui³

(1. College of Information Technology, Shanghai Ocean University, Shanghai 201306, China; 2. Key Laboratory of Fisheries Information, Ministry of Agriculture and Rural Affairs, Shanghai 201306, China; 3. Nantong Longyang Aquatic Products Co., Ltd., Nantong 226634, China)

Abstract: In order to accurately predict the content of dissolved oxygen (DO) in aquaculture water, a prediction model of dissolved oxygen content in aquaculture water based on self-attention mechanism (ATTN) and improved K -means clustering-bidirectional long-term and short-term memory network (BiLSTM) was proposed. Firstly, according to the similarity of environmental data, the improved K -means algorithm was used to divide environmental data into several categories. Then, based on BiLSTM, residual connection was constructed and batch normalization (BN) was added to complete high-level feature extraction, and the feature information was saved by the long-term memory ability of BiLSTM. Finally, the self-attention mechanism was introduced to highlight the importance of data characteristics at different time nodes, which further improved the performance of the model. The experimental results showed that the mean absolute error (MAE), root mean square error

(RMSE) and average absolute percentage error (MAPE) of the hybrid model based on self-attention mechanism and improved K -BiLSTM were 0.238, 0.322 and 0.035, respectively. Compared with single BP model, CNN-LSTM model and traditional K -means-BiLSTM-ATTN model based on residual and BN, the model constructed in this study had better prediction performance and generalization ability.

收稿日期: 2023-01-28

基金项目: 江苏现代农业产业关键技术创新项目[CX(20)2028]; 广东省重点领域研发计划项目(2021B0202070001)

作者简介: 冯国富(1971-), 男, 河南鹤壁人, 博士, 副教授, 研究方向为嵌入式技术研究。(E-mail) gffeng202212@163.com

通讯作者: 陈 明, (E-mail) chengmm202212@163.com

Key words: aquaculture; dissolved oxygen prediction; K-means clustering; bidirectional long-term and short-term memory network (BiLSTM); self-attention mechanism

溶解氧含量是决定水产养殖质量和产量的重要因素。渔业养殖中溶解氧含量过高或不足都会影响鱼类的新陈代谢和繁殖,严重的情况下甚至会影响鱼类的正常生长^[1]。如何利用现有监测数据,准确预测未来溶解氧含量的变化趋势,为养殖人员提供决策参考,已成为近年来国内外学者的研究热点。

由于溶解氧具有时序性、不稳定性和非线性等特点,且受多种因素的影响,各因素之间存在复杂的耦合关系^[2],难以用传统方法^[3-7]和机器学习方法^[8-15]进行建模。而反向传播神经网络(BP)方法又无法有效提取水质气象数据时序维度上的关联,隐藏的时序信息无法被有效利用^[16]。循环神经网络(RNN)适合处理时间序列问题,可有效地关联上下文信息,常用的循环神经网络是长短期记忆网络(LSTM)和门控循环神经网络(GRU),在此基础上,陈英义等^[17]建立了基于 WT-CNN-LSTM 的溶解氧含量预测模型,并取得了不错的效果;曹守启等^[18]提出了改进 LSTM 的水产养殖水体溶解氧含量预测模型并在不同天气条件下成功预测溶解氧含量变化趋势;Wu 等^[19]采用混合 XGBoost-ISSA-LSTM 模型实现对池塘短期和长期溶解氧含量准确预测;Yang 等^[20]构建基于改进鲸鱼优化算法的 GRU 神经网络预测海参养殖水质,提升了模型的精度和泛化能力。但上述模型并没有考虑到相似环境因素下溶解氧含量的变化规律,对于非线性更突出的多参数输入的预测效果并不理想,容易出现预测精度低和结果滞后的情况,因此本研究在基准模型预测之前对输入水质和气象数据进行 K-means 聚类,以更好地反映溶解氧含量变化趋势。

若输入特征之间的关系复杂或者预测时长过长,会导致预测结果滞后和误差增大的问题^[21]。为了筛选与目标关联性更强的参数特征,Yang 等^[22]提出用 CNN-BiLSTM-AM 模型对水产养殖水体溶解氧含量分别进行短期和长期预测;Zhang 等^[23]构建注意力机制和 BiLSTM 模型进行河流水质预测;Li 等^[24]通过引入小波变换和注意力机制提出了一种新的混合模型 BC-MODWT-DA-LSTM 来预测水产养殖水体氨氮含量。上述多阶段模型引入了注意力机制,但因其是通过简单堆叠建立模型,注意力的权重

容易出现偏差,为了防止偏差,本研究采用自注意力机制自适应调整权重,并且在基准模型中加入 BN 层和残差连接防止过拟合。

本研究拟提出一种基于自注意力机制(ATTN)和改进的 K-BiLSTM 的溶解氧含量预测模型,并以江苏省南通市中洋河豚庄园养殖基地的养殖水质和气象数据为样本进行模型训练,提高预测精度。

1 数据获取与预处理

1.1 数据获取

本研究所采用的数据来自南通中洋河豚庄园养殖基地。试验池塘长 15.0 m,宽 15.0 m,水深约 1.2 m。水质数据来自于传感器,空间分辨率为 15.0 m,气象数据来源于气象站。试验采用的是从 2021 年 2 月到 2021 年 7 月传感器记录的水温、pH、溶解氧含量、氨氮含量和盐度(电导率法测定溶解盐质量浓度)等养殖环境数据以及气象站记录的大气温度、湿度、大气压、光照度、风速 5 个气象数据,共 6 940 组数据,采样频率均为 30 min,数据缺失率为 1.17%。短时间内参数浮动范围较小,本研究以 2 h 为时间单位计算各项参数的平均值。

1.2 数据预处理

自动数据采集难免会导致数据缺失,为了保证模型的最终预测效果,需预先对水质数据进行处理。本研究采用线性插值法解决数据缺失的问题。线性插值法的计算公式(1)为:

$$X_k = X_w + \frac{(X_r - X_w)(k - w)}{(r - w)} \quad (1)$$

式中, X_k 为要补齐的缺失值; X_w 为 X_k 前面最近的已知数据; X_r 为 X_k 后面最近的已知数据; w, r 为已知数据时间点, k 为缺失值时间点,位于 w 和 r 之间。

由于采集到的数据各个分量具有不同的维度和量纲,在对数据进行划分数据集之前,首先采用公式(2)对数据进行归一化处理:

$$x = (x' - x_{\min}) / (x_{\max} - x_{\min}) \quad (2)$$

式中, x' 和 x 分别为原始采样数据和归一化之后的数据, x_{\min} 和 x_{\max} 分别为原始数据的最小值和最大值。

2 模型的构建

2.1 改进的 K-means 算法

2.1.1 K-means 算法 K-means 算法是一种迭代聚类算法,它通过将大量样本数据按照到预先选定的聚类中心之间的欧式距离进行聚类,形成多个聚类簇。聚类簇重新选定聚类中心后再次迭代聚类,最终达到最优聚类效果。假设样本集 $X = \{x_1, x_2, x_3, \dots, x_n\}$ 中有 n 个样本,每个样本有 p 个特征参数, $x_i = \{x_{i1}, x_{i2}, x_{i3}, \dots, x_{ip}\}$ 。随机选取 k 个样本, $V = \{v_1, v_2, v_3, \dots, v_k\}$ 作为初始聚类中心,则每个样本 x_i 与初始聚类中心 v_k 的欧式距离为:

$$d_{(x_i, v_k)} = \sqrt{(x_{i1} - v_{k1})^2 + \dots + (x_{ip} - v_{kp})^2} \quad (3)$$

K-means 聚类的目标是使得每个样本到它所在聚类簇的聚类中心的距离尽可能小并且聚类簇间的距离尽可能大,直至达到指定的迭代次数或者聚类中心不再发生变化为止。优化目标为:

$$E = \sum_{i=1}^k \sum_{z_j \in k_i} (z_j - v_k)^2 \quad (4)$$

式中, k_i 为聚类簇样本, z_j 为聚类簇 k_i 中的样本, v_k 为聚类簇 k_i 的簇中心, E 为簇内平方和,其值越小越好。

2.1.2 改进的 K-means 算法

2.1.2.1 轮廓系数的改进 轮廓系数是用于评价聚类算法好坏的一种指标,常用于确定分类簇数值 (k)。轮廓系数包括内聚度和分离度,内聚度代表了类内样本之间的紧密程度,内聚度越低代表聚类结果越紧凑,分离度代表了类间样本的紧密程度,分离度越高代表类间分离程度越高。轮廓系数的计算公式如下:

$$S(i) = \frac{[b(i) - a(i)]}{[a(i), b(i)]_{\max}} \quad (5)$$

其中 $a(i)$ 代表当前样本在同类样本中的平均距离, $b(i)$ 代表当前样本在距离它最近的类别中的平均距离。轮廓系数 $S(i)$ 值越大时 k 值选取得越好。

为了降低传统轮廓系数没有考虑类内的最小距离和类间的平均距离带来的潜在影响,引入了点 i 到它所属类中其他点的最小距离 $s(i)$ 以及点 i 到非所属类中所有点的最大平均距离 $r(i)$,改进后的轮廓系数公式如下:

$$S(i) = \frac{[r(i) - a(i), b(i) - s(i)]_{\max}}{[a(i), b(i)]_{\max}} \quad (6)$$

M 个样本点的平均改进的轮廓系数 S 为:

$$S = \frac{1}{M} \sum_{i=1}^M \frac{[r(i) - a(i), b(i) - s(i)]_{\max}}{[a(i), b(i)]_{\max}} \quad (7)$$

公式(7)反映出各个样本之间更加全面的制约关系。

2.1.2.2 初始聚类中心选取方法的改进 传统聚类方法随机选择 k 个数据作为聚类中心,容易陷入局部最优且无法获得最佳聚类簇。本研究针对此方面进行改进。

数据集中每个样本 x_i 距离其他样本的平均欧式距离是:

$$D(i) = \frac{\sum_{j=1}^N d_{(x_i, x_j)}}{N} \quad (8)$$

其中 x_j 为样本集中的其他样本; N 为样本总数; x_i 为当前样本; $d_{(x_i, x_j)}$ 为样本距离。如果样本 x_i 的平均欧式距离内的样本数量越多,则说明 x_i 是这个样本集所包含的某一区域的中心。此时以 x_i 作为聚类的初始中心更易收敛。计算公式如下:

$$Num[x_i, D(i)] = \sum_{j=1}^N u[d_{(x_i, x_j)} - D(i)] \quad (9)$$

式中, $Num[x_i, D(i)]$ 为以 x_i 为中心,以 $D(i)$ 为半径范围内样本点的个数; $d_{(x_i, x_j)}$ 为当前样本距离聚类中心的距离; $u[d_{(x_i, x_j)} - D(i)]$ 为阶跃函数, $u(v)$ 公式为:

$$u(v) = \begin{cases} 1, & v \geq 0 \\ 0, & v < 0 \end{cases} \quad (10)$$

v 为函数自变量。计算出所有样本 x_i 的 $Num[x_i, D(i)]$ 值之后,从大到小选取前 k 个作为初始聚类簇中心,可有效避免陷入局部最优。

2.1.2.3 相似度量方法改进 在 K-means 聚类分析中,可以使用欧式距离和余弦距离来衡量水质参数 X 和任意其他水质数据 Y 之间的相似度,这两者都包含 m 维特征,如水温、氨氮含量。欧式距离用来衡量多维空间中点之间的实际绝对距离,反映样本属性之间的数值差异性。公式如下:

$$d_{(x, y)} = \sqrt{\sum_{j=1}^m (x_j - y_j)^2} \quad (11)$$

式中 x_j 和 y_j 分别为 X 和 Y 的第 j 个分量。余弦距离指的是,计算 2 点与原点所成的直线夹角的余弦值,取值范围是 $[-1, +1]$,越趋近于 1 代表越相似,越趋近于 -1 代表方向相反,0 代表正交,简单来说,夹角越大就说明两点越不相似,夹角越小说明越相似。其公式如下:

$$\cos(X, Y) = \frac{\sum_{j=1}^m x_j y_j}{\sqrt{\sum_{j=1}^m x_j^2 \sum_{j=1}^m y_j^2}} \quad (12)$$

鉴于 2 种度量方式优势互补,提出改进的相似度量公式:

$$\text{sim}(X, Y) = \frac{1}{1+d_{(X,Y)}} \cdot \cos(X, Y) \quad (13)$$

改进的相似度综合考虑了欧式距离相似度和余弦相似度,且存在上限。 $\text{sim}(X, Y)$ 值越大说明 2 点越相似, $\text{sim}(X, Y)$ 值越小说明越不相似。

2.1.2.4 改进的 K-means 算法流程 为了避免不良样本的引入导致模型的收敛速度慢、预测精度低等问题,本研究首先对归一化之后的数据采用改进后的 K-means 聚类^[25]将相似的样本划分成一簇,然后在相同的簇中建立模型来进行溶解氧含量的预测。对于给定的 n 个样本所构成的集合,选取改进的轮廓系数最大值所对应的 k 值作为样本簇的个数,按照改进的初始聚类中心选取方法来选取初始聚类中心,根据当前样本距离簇中心的综合相似度距离,将其划分到最近的簇,然后一直迭代至簇内的样本距离簇中心的综合距离尽可能小,相似度尽可能高,簇间的距离尽可能大,相似度尽可能低。算法流程如下:

- 1) 为 n 个样本计算 $\text{Num}[x_i, D(i)]$ 值并从大到小排序,选取前 k 个点作为初始聚类中心,即簇中心。
- 2) 计算每个样本到各个簇中心的综合相似度距离,将其划分至综合距离最近的簇。
- 3) 更新每个簇中心。
- 4) 重复步骤 2~步骤 3 直至每个簇中心不再发生变化。归一化之后的样本经过改进的 K-means 聚类之后得到 k 个簇,按照公式(14)计算出每个簇的中心。

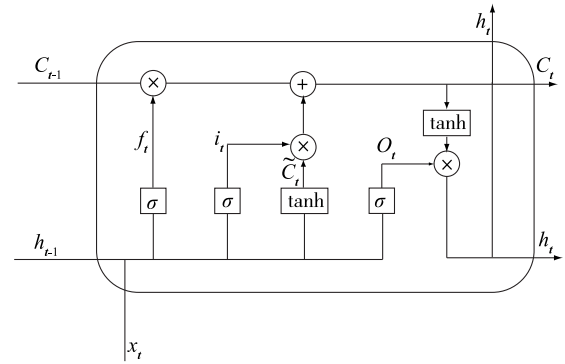
$$\bar{x}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij} \quad (14)$$

其中第 i 簇含 n_i 个样本, x_{ij} 表示第 i 类的第 j 个样本。模型测试过程中,计算当前测试集中的样本到各个簇中心的综合相似度距离,将距离最近的簇作为当前样本的所属簇,使用所属簇的神经网络模型对当前样本进行溶解氧含量预测。

2.2 改进的双向长短期记忆网络 (BiLSTM)

2.2.1 长短期记忆网络 (LSTM)

LSTM 神经网络可以有效地捕捉到长期顺序特征,而无需进行特征工程,此外,它解决了 RNN 存在的梯度消失和短期记忆问题,使模型能够正常收敛。因此, LSTM 已经成为时间序列建模的一种重要工具。LSTM 的基本结构如图 1 所示。LSTM 拥有 3 种类型的门结构(遗忘门、输入门和输出门)来保护和控制细胞状态。



x_t 表示当前时刻的输入; h_{t-1} 、 C_{t-1} 分别表示上一个时刻的输出和细胞状态; f_t 、 i_t 、 O_t 分别为遗忘门、输入门和输出门的输出; \tilde{C}_t 为候选细胞状态; h_t 和 C_t 分别为当前时刻的输出和细胞状态; σ 和 \tanh 为激活函数。

图 1 长短期记忆网络 (LSTM) 的结构

Fig.1 Structure diagram of long-term and short-term memory network (LSTM)

2.2.1.1 遗忘门 在 LSTM 中, 上一个输出信号会通过遗忘门, 该门决定会从细胞状态中保留和舍弃什么信息。遗忘门根据上一个输出 h_{t-1} 和当前输入 x_t 进行 Sigmoid 非线性映射, 并输出一个各个分量都在 0 到 1 之间的向量 f_t , 1 表示完全保留, 0 表示完全舍弃, 最后与细胞状态 C_{t-1} 相乘。其公式如下所示:

$$f_t = \sigma[W_f \cdot (h_{t-1}, x_t) + b_f] \quad (15)$$

式中, x_t 为当前时刻输入; h_{t-1} 为上一时刻输出; W_f 和 b_f 为权重和偏置; f_t 为遗忘门输出; σ 为激活函数。

2.2.1.2 输入门 LSTM 的输入门是一种用于控制信息流的门, 它由一个 Sigmoid 函数和一个乘法运算组成, 通过对当前输入和上一时刻状态隐藏向量的组合来控制信息的流动。输入门控制哪些信息能够进入 LSTM 单元, 而哪些信息不能进入, 从而控制 LSTM 单元的输出, 其公式如公式(16)和公式(17)所示:

$$i_t = \sigma[W_i \cdot (h_{t-1}, x_t) + b_i] \quad (16)$$

$$\tilde{C}_t = \tanh[W_c \cdot (h_{t-1}, x_t) + b_c] \quad (17)$$

式中, x_t 为当前时刻输入; h_{t-1} 为上一时刻输出;

W_i 、 W_c 为权重; b_i 和 b_c 为偏置; i_t 为输入门的输出; \tilde{C}_t 为候选细胞状态。

2.2.1.3 记忆单元更新 将旧细胞状态 C_{t-1} 更新为 C_t , 其更新公式为:

$$C_t = C_{t-1} \cdot f_t + \tilde{C}_t \cdot i_t \quad (18)$$

式中, C_t 为当前时刻细胞状态; C_{t-1} 为上一时刻细胞状态; f_t 为遗忘门输出; \tilde{C}_t 为候选细胞状态; i_t 为输入门的输出。

2.2.1.4 输出门 LSTM 输出门也是一种用于控制信息流的门, 它由一个 Sigmoid 函数和一个乘法运算组成。它控制什么时候将 LSTM 单元当前的输出信息发送出去, 以及何时将该信息保留在单元内部。输出门也可以控制 LSTM 单元的输出, 使得输出信息的流动更加有序, 从而达到更好的预测结果, 其公式如公式 (19) 和 公式 (20) 所示:

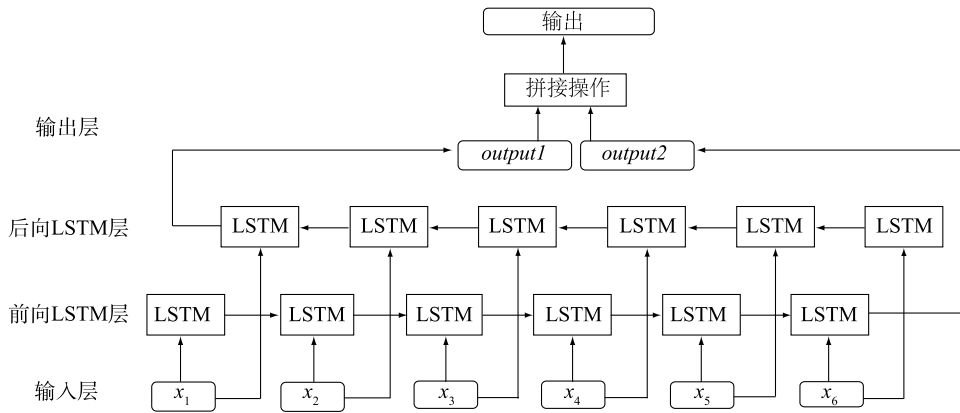
$$o_t = \sigma[W_o \cdot (h_{t-1}, x_t) + b_o] \quad (19)$$

$$h_t = o_t \cdot \tanh(C_t) \quad (20)$$

式中, x_t 为当前时刻输入; h_{t-1} 为上一时刻输出; W_o 、 b_o 为权重和偏置; C_t 当前时刻细胞状态; o_t 为输出门的输出; h_t 为当前时刻输出。

输出门将内部状态的信息传递给外部状态 h_t , Sigmoid 函数确定记忆单元的哪些信息传递出去, 然后, 细胞状态通过 \tanh 函数得到 $[-1, 1]$ 的值并将它和输出门的输出相乘, 最终外部状态仅仅会得到输出门确定输出的那部分。

2.2.2 双向长短期记忆网络 (BiLSTM) 双向长短期记忆神经网络是一种用于处理和预测时间序列数据的深度学习架构, 它开发自传统的长短期记忆网络 (LSTM)。BiLSTM 的基本思想是, 它使用 2 个独立的隐藏层, 分别处理前向和后向数据流。这样, BiLSTM 可以更好地捕获输入数据中的历史和未来信息, 从而更好地处理时间序列数据。其基本流程如图 2 所示。



x_t 表示当前时刻的输入; $output1$ 、 $output2$ 分别表示后向 LSTM 输出和前向 LSTM 输出。LSTM: 长短期记忆网络。

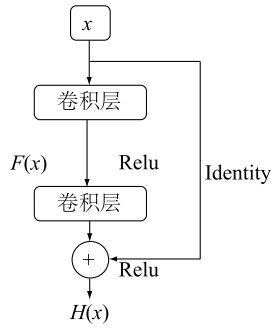
图 2 双向长短期记忆网 (BiLSTM) 的结构

Fig.2 Structure diagram of bidirectional long-term and short-term memory network (BiLSTM)

2.2.3 BiLSTM 的改进 残差连接是一种对神经网络进行深度构建的技术, 它可以将多个神经网络层连接起来, 并通过跳过连接来减少模型的复杂性, 从而改善网络的性能。残差连接的思想是: 如果模型中的某些层不起作用, 那么它们可以被简单地“跳过”, 从而帮助模型学习更深层次的特征, 有效地减少训练时间和提高准确率。本研究以构建残差块 (Residual block) 的方式大大加深 BiLSTM 模型层次。基本残差结构如图 3 所示。

随着神经网络隐藏层的加深, 容易产生梯度消

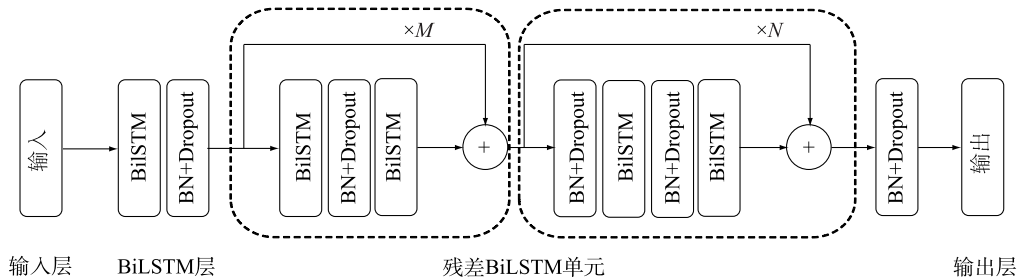
失和梯度弥散的问题, 而通过残差连接可以有效避免这 2 种情况的出现, 并且能够减少网络模型的参数量, 提高模型的训练效率, 有助于高效提取更高层次的特征。BiLSTM 能够有效提取时序数据特征, 本研究结合批标准化 (BN) 和 BiLSTM 构建出基于残差和引入 BN 层的 BiLSTM 网络, 模型如图 4 所示。模型包括 BiLSTM 层和多个残差单元, 其共同完成溶解氧含量相关环境数据特征提取, 随后将特征送入 BN 和随机丢弃 (Dropout), 进一步提高模型训练效率和增强泛化能力。



x 表示输入; $F(x)$ 和 $H(x)$ 分别为神经网络块的映射函数和输出; ReLU 为激活函数。Identity 表示输入 x 经过变换后与卷积层的输出结果相加。

图 3 基本残差结构

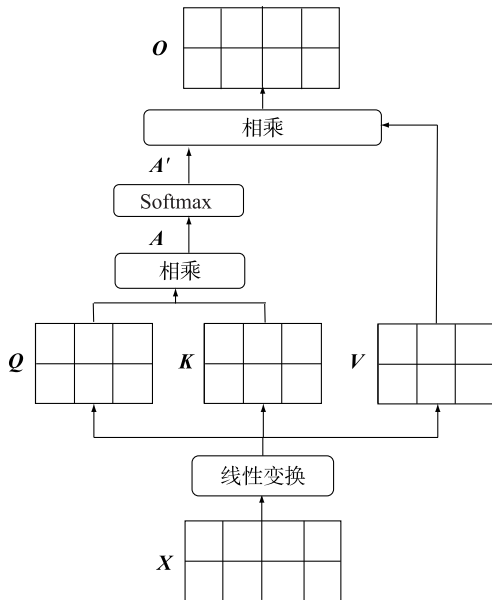
Fig.3 Basic residual structure



Dropout: 随机丢弃; BiLSTM 为双向长短期记忆网络模块; BN 为批标准化。 M 、 N 表示基准模型的堆叠个数。

图 4 基于残差和引入批标准化 (BN) 的 BiLSTM 网络结构

Fig.4 BiLSTM network structure diagram based on residual and batch normalization (BN)



X 表示输入矩阵; Q 、 K 、 V 分别表示查询矩阵、键矩阵和值矩阵; A 和 A' 为中间层输出; O 为输出矩阵; Softmax 为激活函数。

图 5 自注意力机制的基本结构

Fig.5 Basic structure diagram of self-attention mechanism

2.3 自注意力机制

注意力机制 (Attention mechanisms) 技术可以追溯到人类视觉系统。当人们观察时, 会有选择性地捕获重要的信息, 而忽略不太重要的内容^[26]。注意力机制旨在通过将有限的资源集中在处理更重要的信息上, 以提升神经网络在不同时间点上对重要信息的关注能力, 进而提高模型的性能。自注意力机制是一种利用数据特征内在信息的注意力机制, 以便在多个时间节点上发现输入特征之间的相关性。其基本结构如图 5 所示。

首先, 自注意力机制输入矩阵 X 经过线性变换得到矩阵 Q 、 K 、 V , 变换矩阵是通过学习得到的。其次, 矩阵 Q 和矩阵 K 的转置相乘, 再除以 1 个尺度 $\sqrt{d_k}$ 得到相关性矩阵 A , 其中, d_k 为查询向量和键向量之间的维度, 目的是为了防止方差过大, 分布陡峭。然后使用 Softmax 函数将矩阵 A 归一化为 A' 。最后将 A' 与矩阵 V 相乘得到注意力机制层的输出特征。计算公式如下:

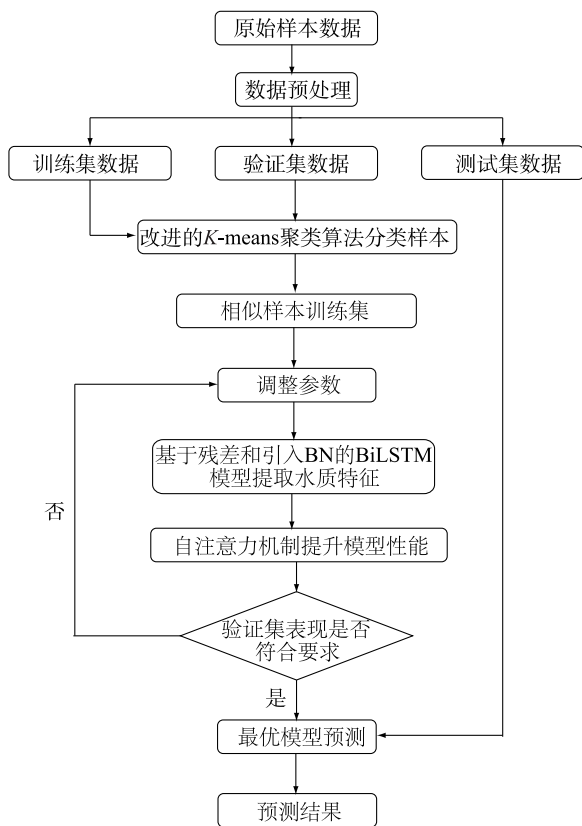
$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (21)$$

式中, Q 、 K 、 V 分别是查询矩阵、键矩阵和值矩阵; $softmax$ 为激活函数; $\sqrt{d_k}$ 为缩放因子, d_k 是 K 的维度; $Attention(Q, K, V)$ 为模型输出。本研究使用自注意力机制的目的是为了获得时间维度上水质特征的重要程度, 以提高模型的预测性能。

2.4 预测模型

养殖过程中水体溶解氧含量的变化受多个复杂因素影响, 并且在时间维度上呈非线性变化趋势, 所以在溶解氧含量预测任务中, 不仅要考虑如水温、

pH 等参数之间的复杂关系,而且模型构建还要考虑到时间维度上的变化。BiLSTM 能够有效地捕捉长序列之间的数据关联,缓解梯度消失现象,同时效果优于传统 RNN,比 LSTM 能够更有效提取时序数据特征,因此本研究使用 BiLSTM 作为基准模型。为了有效避免梯度消失和提取高层次特征以及增强模型的泛化能力,本研究在 BiLSTM 的基础上引入了残差和 BN 层。由于溶解氧含量在时间维度上呈现差异性,在上述改进的 BiLSTM 模型基础之上加入自注意力机制 ATTN,获取不同时间节点上的特征重要程度,降低非重要特征的影响,提高预测精度。预测流程如图 6 所示。



BiLSTM 为双向长短期记忆网络模块;BN 为批标准化。

图 6 改进的 K-BiLSTM-ATTN 模型的流程

Fig.6 Flow chart of improved K-BiLSTM-ATTN model

使用本研究模型对溶解氧含量预测的步骤如下:

1) 将获取的水质数据进行缺失值填充和归一化处理,按照 6:2:2 的比例划分成训练集、验证集和测试集。

2) 对训练集和验证集使用改进的 K-means 算法划分类别,对不同类别的数据集分别进行模型构建和训练。

3) 初始化基于残差和 BN 层的 BiLSTM-ATTN 模型参数,将训练集数据输入模型不断训练直至验证集表现良好或达到预定训练轮次。

4) 将测试集中的数据输入所属类别模型中得到测试集的预测结果,将本研究模型与其他模型的预测结果进行对比分析,得到溶解氧含量最优预测模型。

2.5 模型评价指标

本研究选取平均绝对误差 (MAE)、均方根误差 (RMSE)、平均绝对百分比误差 (MAPE) 作为模型的性能评估指标,其计算公式为公式 (22) ~ 公式 (24)。

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \bar{y}_i| \quad (22)$$

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \bar{y}_i)^2} \quad (23)$$

$$MAPE = \frac{1}{N} \sum_{i=1}^N \frac{|y_i - \bar{y}_i|}{y_i} \quad (24)$$

式中, N 表示样本个数, y_i 和 \bar{y}_i 分别表示预测值和真实值。以上 3 个指标的取值范围是 $[0, \infty)$, $RMSE$ 值越小表明模型的预测能力越好, MAE 和 $MAPE$ 值越小表明模型稳定性越好。

3 结果与分析

3.1 改进的 K-means 聚类算法评估

我们采用改进的轮廓系数来选取聚类簇的个数 (k), k 取不同值时的轮廓系数值如表 1 所示。当 k 取 4 时轮廓系数最大,因此本研究采用 4 作为聚类簇的个数。

表 1 不同簇数的轮廓系数值

Table 1 The values of the contour coefficient for different cluster numbers

k	2	3	4	5	6	7	8	9	10
S	0.487	0.506	0.508	0.465	0.467	0.472	0.468	0.462	0.467

k : 聚类簇的个数; S : 轮廓系数值。

3.2 模型参数

所有试验均在 PC 主机上运行,主机性能: 2.3 GHz Intel i7-11800H 处理器、16 G 内存、NVIDIA Ge-

Force RTX3060 显卡、Microsoft Windows 10。选用 Tensorflow 作为深度学习平台。本研究使用 2.1 节提出的改进的 K -means 聚类算法将训练集和验证集数据划分成 4 个簇,并分别在各个簇中训练基于残差和 BN 的 BiLSTM-ATTN 模型。溶解氧含量预测模型主要由改进的 K -menas 聚类、BiLSTM 层、残差、BN 层、ATTN 层、全连接层和输出层构成。模型结构确定后,多次调整模型参数进行试验,得到最优的基于自注意力机制和改进的 K -BiLSTM 模型。表 2 是本研究模型的参数设置。

表 2 模型主要参数设置

Table 2 The main parameters of the model

模型参数	设置值(函数)
BiLSTM 层神经元	256 个
BiLSTM 激活层函数	Relu
Self-Attention 激活函数	Sigmoid
Dense 层神经元	32 个
所有子模型卷积核大小	3×3
批大小	128 个样本
迭代次数	2 000 次
学习率	0.000 2

模型优化器为 Adam。

3.3 模型对比分析

为验证本研究模型的性能优势,将本研究模型与单一的 BP 模型、CNN-LSTM 模型、传统的 K -means-基于残差和 BN 的 BiLSTM-ATTN 模型、改进的 K -means-BiLSTM-ATTN 模型和未引入自注意力机制的改进的 K -BiLSTM 5 个模型进行对比。使用 MAE 、 $RMSE$ 和 $MAPE$ 统计指标对不同模型进行评估和比较,结果如表 3 所示。

由表 3 中 Model1 和 Model2 与 Model6 的对比可以得出,预先对样本进行相似聚类可以有效剔除具有较大差异的样本数据进而提升模型的预测精度。其次,Model3 和 Model6 相比可以得出,改进的 K -means 聚类具有更好的聚类效果并且对各个类别的模型预测能力提升贡献较大。由 Model4 和 Model6 相比的结果可以得出结论,残差连接的构建和 BN 层的加入使得本研究模型拥有更强的特征提取能力和泛化能力。比较 Model5 和 Model6 模型评价指标,表明自注意力机制的引入提升了模型的预测能力和稳定性。本研究所提出的改进的 K -means-基于

残差和 BN 的溶解氧含量预测模型的 MAE 、 $RMSE$ 、 $MAPE$ 分别为 0.238、0.322 和 0.035,表明,本研究提出的混合模型在溶解氧含量预测方面优于单一的 BP 模型、CNN-LSTM 模型、传统的 K -means-基于残差和 BN 的 BiLSTM-ATTN 模型、改进的 K -means-BiLSTM-ATTN 模型和未引入自注意力机制的改进的 K -BiLSTM 模型。

表 3 6 种不同预测模型评价指标对比

Table 3 Comparison of evaluation indices of six different prediction models

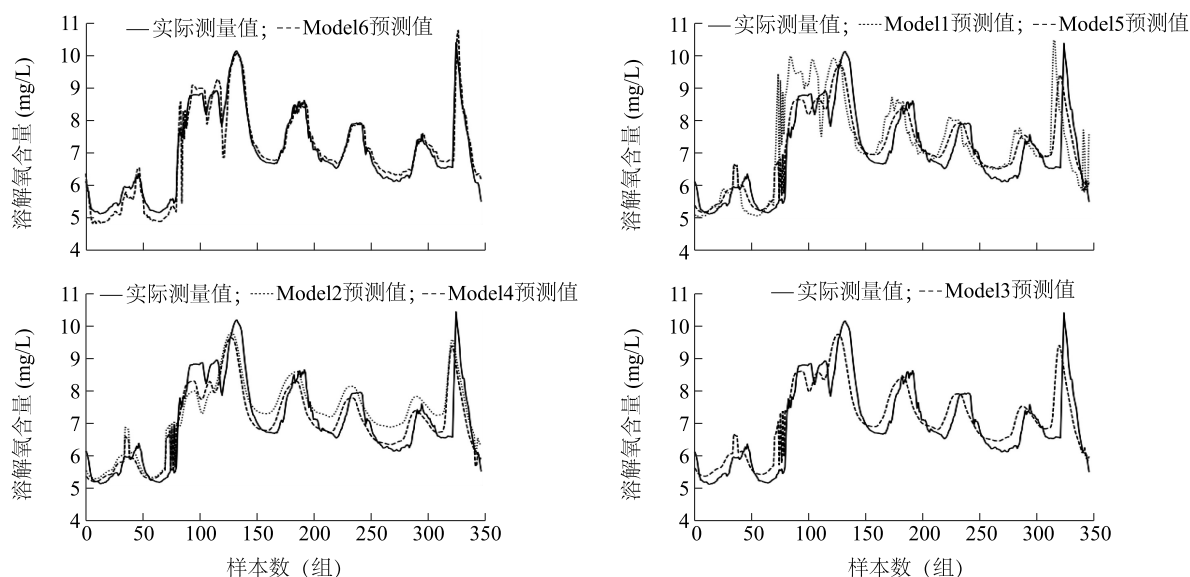
模型	MAE	$RMSE$	$MAPE$
Model1	0.806	1.134	0.114
Model2	0.558	0.705	0.080
Model3	0.467	0.658	0.067
Model4	0.416	0.603	0.058
Model5	0.416	0.589	0.060
Model6	0.238	0.322	0.035

Model1~Model6 分别为单一的 BP 模型、CNN-LSTM 模型、传统的 K -means-基于残差和 BN 的 BiLSTM-ATTN 模型、改进的 K -means-BiLSTM-ATTN 模型、未引入自注意力机制的改进的 K -BiLSTM 模型和本研究提出的改进的 K -means-基于残差和 BN 的 BiLSTM-ATTN 模型。 MAE 为平均绝对误差; $RMSE$ 为均方根误差; $MAPE$ 为平均绝对百分比误差。

图 7 为各预测模型溶解氧含量的预测结果和实际溶解氧含量的比较,其中横坐标是样本数,共 347 组测试数据,纵轴为溶解氧含量(mg/L)。从图中可以看出,相比于其他模型,本研究提出的改进的 K -means-基于残差和 BN 的 BiLSTM-ATTN 模型的预测结果波动更小更贴近真实曲线。

4 结论

为了提升水产养殖过程中水体溶解氧含量的预测精度,本研究提出了一种基于自注意力机制和改进的 K -BiLSTM 的溶解氧含量预测的混合模型。采用改进的 K -means 聚类对样本数据进行分类,防止不同类别样本间的过度干扰,提高了预测数据源的准确性。引入了残差连接和 BN 层,不仅有效地减少训练时间和提高准确率,还起到提取更高层次特征的作用。BiLSTM 网络和自注意力机制分别起到了在时间序列上进行长期记忆保存与突出重要信息的作用。本研究提出的模型具有更高的精确度和较好的鲁棒性,可用于实际渔业生产。



图中 Model1~Model6 分别为单一的 BP 模型、CNN-LSTM 模型、传统的 K -means-基于残差和 BN 的 BiLSTM-ATTN 模型、改进的 K -means-BiLSTM-ATTN 模型、未引入自注意力机制的改进的 K -BiLSTM 模型和本研究提出的改进的 K -means-基于残差和 BN 的 BiLSTM-ATTN 模型。

图7 溶解氧含量预测结果与实际结果对比图

Fig.7 Comparison of dissolved oxygen content between predicted results and actual results

参考文献:

- [1] LIPIZER M, PARTESCANO E, RABITTI A, et al. Qualified temperature, salinity and dissolved oxygen climatologies in a changing Adriatic Sea[J]. Ocean Science, 2014, 10(5): 771-797.
- [2] 陈英义, 程倩倩, 方晓敏, 等. 主成分分析和长短时记忆神经网络预测水产养殖水体溶解氧[J]. 农业工程学报, 2018, 34(17): 183-191.
- [3] 金光炎. 水文统计理论与实践[M]. 南京: 东南大学出版社, 2012.
- [4] 刘明, 李由明, 王平, 等. 基于小波分解的凡纳滨对虾养殖水体水质的仿真研究[J]. 广东农业科学, 2013, 40(17): 170-172.
- [5] 徐梅, 晏福, 刘振忠, 等. 灰色 GM(1,1)-小波变换-GARCH 组合模型预测松花江流域水质[J]. 农业工程学报, 2016, 32(10): 137-142.
- [6] 岳遥, 李天宏. 基于模糊集理论的马尔可夫模型在水质定量预测中的应用[J]. 应用基础与工程科学学报, 2011, 19(2): 231-242.
- [7] 黄廷林, 韩晓刚, 卢金锁. 基于 Lyapunov 指数的混沌预测方法及在水质预测中的应用[J]. 西安建筑科技大学学报(自然科学版), 2008, 40(6): 846-851.
- [8] ALVAREZ MEZA A M, DAZA SANTACOLOMA G. Parameter selection in least squares support vector machines regression oriented, using generalized cross-validation[J]. Dyna-Colombia, 2012, 79(171): 23-30.
- [9] 邹志红, 王学良. BP 模型在河流水质预测中的误差分析[J]. 环境科学学报, 2007, 27(6): 1038-1042.
- [10] AMID S, GUNDOSHMIAN T M. Prediction of output energies for broiler production using linear regression, ANN (MLP, RBF), and ANFIS models[J]. Environmental Progress & Sustainable Energy, 2017, 36(2): 577-585.
- [11] 刘东君, 邹志红. 最优加权组合预测法在水质预测中的应用研究[J]. 环境科学学报, 2012, 32(12): 3128-3132.
- [12] 刘双印, 徐龙琴, 李振波, 等. 基于 PCA-MCAFA-LSSVM 的养殖水质 pH 值预测模型[J]. 农业机械学报, 2014, 45(5): 239-246.
- [13] 龚怀瑾, 毛力, 杨弘. 基于变尺度混沌 QPSO-LSSVM 的水质溶氧预测建模[J]. 计算机与应用化学, 2013, 30(3): 315-318.
- [14] 孙伯寅, 董国庆, 张荣. 支持向量机在水源水化学耗氧量预测中的应用[J]. 环境与健康杂志, 2016, 33(6): 544-547.
- [15] 罗华军, 黄应平, 刘德富. 基于 WA-SVM 的水库溶解氧预测[J]. 西北农林科技大学学报(自然科学版), 2009, 37(3): 181-186.
- [16] 宦娟, 刘星桥. 基于 K -means 聚类和 ELM 神经网络的养殖水质溶解氧预测[J]. 农业工程学报, 2016, 32(17): 174-181.
- [17] 陈英义, 方晓敏, 梅思远, 等. 基于 WT-CNN-LSTM 的溶解氧含量预测模型[J]. 农业机械学报, 2020, 51(10): 284-291.
- [18] 曹守启, 周礼馨, 张铮. 采用改进长短时记忆神经网络的水产养殖溶解氧预测模型[J]. 农业工程学报, 2021, 37(14): 235-242.
- [19] WU Y H, SUN L Q, SUN X B, et al. A hybrid XGBoost-SSA-LSTM model for accurate short-term and long-term dissolved oxy-

- gen prediction in ponds[J]. Environ Sci Pollut Res Int, 2021, 29 (12): 18142-18159.
- [20] YANG H H, LIU S E. Water quality prediction in sea cucumber farming based on a GRU neural network optimized by an improved whale optimization algorithm [J]. PeerJ Comput Sci, 2022, 8; e1000.
- [21] ZOU Q H, XIONG Q Y, LI Q D, et al. A water quality prediction method based on the multi-time scale bidirectional long short-term memory network [J]. Environmental Science and Pollution Research, 2020, 27(9): 16853-16864.
- [22] YANG W B, LIU W, GAO Q. Prediction of dissolved oxygen concentration in aquaculture based on attention mechanism and combined neural network [J]. Math Biosci Eng, 2023, 20(1): 998-1017.
- [23] ZHANG Q, WANG R Q, QI Y, et al. A watershed water quality prediction model based on attention mechanism and Bi-LSTM [J]. Environmental Science and Pollution Research, 2022, 29(50): 75664-75680.
- [24] LI Y T, LI R. Predicting ammonia nitrogen in surface water by a new attention-based deep learning hybrid model [J]. Environmental Research, 2023, 216: 114723.
- [25] CAO X K, LIU Y R, WANG J P, et al. Prediction of dissolved oxygen in pond culture water based on K -means clustering and gated recurrent unit neural network [J]. Aquacultural Engineering, 2020, 91: 102122.
- [26] 何津民, 张丽珍. 基于自注意力机制和 CNN-LSTM 深度学习的对虾投饵量预测模型 [J]. 大连海洋大学学报, 2022, 37(2): 304-311.

(责任编辑:陈海霞)