

PRML第三次作业报告：Air Pollution Forecasting - LSTM Multivariate

Ziheng Wang
wangziheng@buaa.edu.cn

Abstract

本研究基于Kaggle平台提供的2010年1月2日至2014年12月31日北京地区逐小时空气质量与气象观测数据，构建了一种多变量长短期记忆网络（Long Short-Term Memory, LSTM）模型，用于预测未来一小时的PM2.5浓度。数据集涵盖18项特征，包括PM2.5浓度、气温、露点温度、气压、风向、风速、降水时长等，通过数据清洗、缺失值处理、分类变量数值化以及标准化等步骤完成预处理。在特征工程中，引入滞后11小时的多维气象变量作为输入，目标为当前时刻的PM2.5浓度，从而构建适用于多步预测的时间序列模型。模型采用双层LSTM架构，融合Dropout抑制过拟合，借助Adam优化器和均方误差（MSE）进行训练，并通过早停机制和模型检查点优化训练流程。

相比传统循环神经网络（RNN），LSTM通过引入门控机制和细胞状态，有效解决梯度消失与爆炸问题，具备建模长期依赖关系的能力。实验结果显示，在测试集上，该模型取得了较低的均方误差和平均绝对误差，能够准确跟踪PM2.5浓度的时序变化，尤其在极端天气条件下的预测效果显著优于传统线性方法。可视化分析进一步验证了模型对气象因素驱动下的污染趋势响应能力。本研究不仅展示了LSTM在多变量时间序列预测中的有效性，也为城市空气污染监测与预警提供了可靠的技术支撑，具有良好的应用前景。

Introduction

随着城市化与工业化进程的不断加速，空气污染已成为全球尤其是发展中国家亟待解决的重大环境问题之一。其中，细颗粒物（PM2.5）因其对人类呼吸系统和心血管健康造成显著威胁，引发了广泛关注。以北京为例，PM2.5浓度在多个年份长期维持在较高水平，严重影响居民生活质量与公共健康。因此，构建一种能够实现小时级预测的高精度模型，

对城市空气污染预警与治理具有重要意义。然而，由于污染物浓度受多种气象变量的综合作用，如温度、湿度、气压、风向和降水等，传统基于单变量或线性假设的统计模型难以准确刻画其中的非线性动态演化过程。

在时间序列建模领域，长短期记忆网络（Long Short-Term Memory, LSTM）作为循环神经网络（RNN）的一种变体，因其独特的门控机制与长期依赖建模能力，在自然语言处理、金融市场预测和气象模拟等任务中取得了优异表现。LSTM能够有效缓解长序列训练中的梯度消失问题，适合处理高维、多变量的时序数据。基于这一优势，本文提出采用多变量LSTM网络结构，对北京地区五年逐小时记录的空气质量与气象数据进行建模，利用前10小时的历史记录预测未来一小时的PM2.5浓度。研究旨在挖掘气象与污染之间的隐含关系，提升预测精度，为政府管理与公众防护提供科学决策支持，同时也为复杂环境下的时间序列建模提供一种可行的深度学习路径。

Methodology

本部分将详细阐述问题解决过程中所用到的基于长短期记忆网络（LSTM）架构模型。

M1: The Vanishing/Exploding Gradient Problem

在循环神经网络（Recurrent Neural Network, RNN）中，梯度困境主要表现为梯度消失（vanishing gradient）与梯度爆炸（exploding gradient）问题，是制约其在长序列学习中性能的关键障碍。RNN通过时间步间的参数共享，在处理序列数据时具备捕捉时间依赖关系的能力。然而，这种结构在反向传播过程中需对损失函数相对于每一时间步的参数进行梯度计算，导致误差信号在网络中反复传递时发生指数级衰减或增长。

梯度消失问题指的是，当时间序列较长时，反向传播过程中的梯度值逐渐趋近于零，导致网络难以对早期输入或长期依赖关系进行有效学习。反之，若权重反复乘积的结果大于1，梯度则可能呈指数级增长，导致梯度爆炸，从而引发模型参数剧烈震荡或训练不收敛。这些问题不仅影响模型的收敛速度与稳定性，还限制了RNN在实际应用中的可扩展性，尤其在语言建模、时间序列预测等需要跨时间跨度建模的任务中尤为突出。

为缓解这一困境，研究者提出了多种改进方案，其中最具代表性的便是长短期记忆网络（Long Short-Term Memory, LSTM）和门控循环单元（Gated Recurrent Unit, GRU）。这类模型通过引入门控结构，有效控制信息流动与遗忘机制，从而在较长序列中保持稳定的梯度传递，显著改善了传统RNN在处理长期依赖问题时的性能。

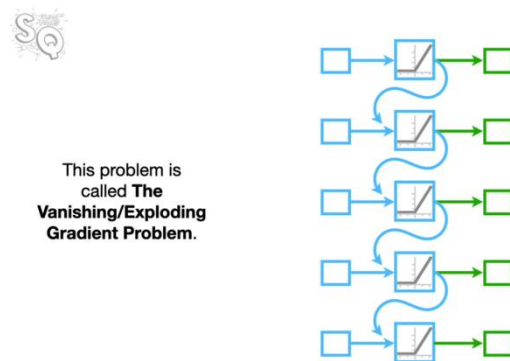


图 1 梯度消失与梯度爆炸（图片来源：StatQuest with Josh Starmer）

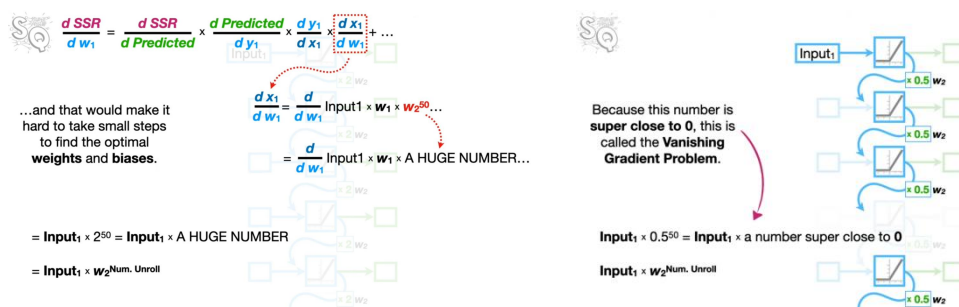


图 2 梯度消失与梯度爆炸数学分析（图片来源：StatQuest with Josh Starmer）

M2: LSTM Model

长短期记忆网络（Long Short-Term Memory, LSTM）作为一种结构优化的循环神经网络，其核心优势在于通过门控机制调控信息的选择性传递，从而有效应对传统RNN在长序列训练中所遇到的梯度消失与梯度爆炸问题。LSTM单元的数学结构围绕四个关键组件展开：遗忘门（forget gate）、输入门（input gate）、候选状态（candidate state）以及输出门（output gate）。

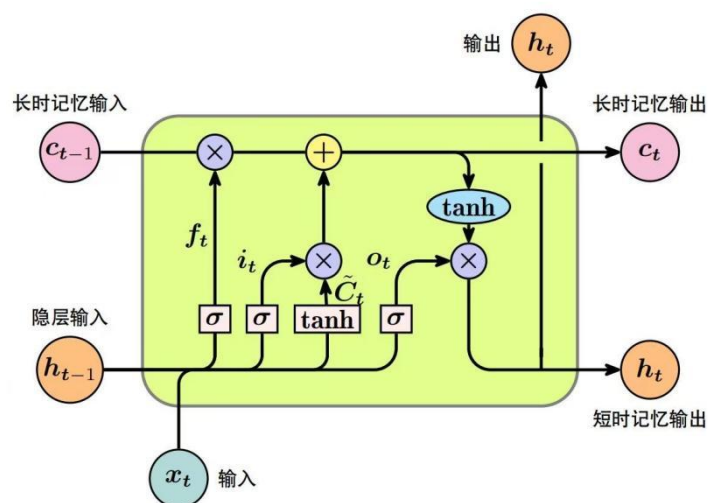


图 3 LSTM Model

1、遗忘门（Forget Gate）

遗忘门负责控制上一时刻的记忆状态 C_{t-1} 在当前时刻是否被保留或遗弃。它的输出是一个值域在 $[0,1]$ 之间的向量，表示每一维记忆内容的保留比例。其计算表达式为：

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

其中， x_t 为当前输入， h_{t-1} 为上一时刻的隐藏状态， W_f 和 b_f 分别是权重矩阵与偏置项， σ 为sigmoid激活函数。输出 f_t 接近1表示信息被保留，接近0表示信息被遗忘。

2、输入门与候选状态（Input Gate and Candidate State）

输入门决定当前输入信息写入到记忆单元的程度，配合候选状态共同更新记忆内容。其数学形式如下：

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

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

其中， i_t 是输入门的激活结果，决定候选状态 \tilde{C}_t 中的信息有多少被引入当前的细胞状态。候选状态通过 \tanh 函数生成，提供一种新的信息候选，用于更新记忆。

3、更新记忆单元状态

在获得遗忘门与输入门的输出后，LSTM更新其细胞状态 C_t ，融合保留的历史信息与当前新的输入信息：

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$$

其中， \odot 表示按元素乘法。这一结构使得信息可以在时间步之间长期存储，且不会被梯度过度缩放或爆炸。

4、输出门（Output Gate）

输出门决定当前记忆状态中哪些部分用于生成新的隐藏状态 h_t ，从而影响最终的输出：

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

$$h_t = o_t \odot \tanh(C_t)$$

输出门的输出 o_t 调控着细胞状态 C_t 经过非线性激活函数后的输出量，进而决定当前时刻的隐藏表示。

M3: Application in Air Pollution Prediction

牛顿法是一种二阶优化方法，前面的模型假设和损失函数与前面均相同，只是参数更新公式不同。其利用梯度和Hessian矩阵（二阶导数矩阵）来更新参数，更新公式为：

本研究基于多变量时间序列建模思路，引入LSTM神经网络对PM2.5浓度进行预测，模型输入包含细颗粒物历史数据及多种气象参数（如温度、湿度、风速等）。为增强模型对多尺度时序特征的提取能力，我构建了一个由两层LSTM组成的深层网络架构，分别包含32个与16个单元，用以逐层提炼污染与气象变量之间的非线性依赖关系。最终，经过全连接层输出预测结果，实现对未来污染水平的精准估计。

Model: "sequential"

Layer (type)	Output Shape	Param #
lstm (LSTM)	(None, 11, 32)	5,120
dropout (Dropout)	(None, 11, 32)	0
lstm_1 (LSTM)	(None, 16)	3,136
dense (Dense)	(None, 1)	17

Total params: 8,273 (32.32 KB)
Trainable params: 8,273 (32.32 KB)
Non-trainable params: 0 (0.00 B)

图 4 模型架构

Experimental Studies

E1: Data Cleaning and Preprocessing

原始数据如下：

information of train set:

	date	pollution	dew	temp	...	wnd_dir	wnd_spd	snow	rain
0	2010-01-02 00:00:00	129.0	-16	-4.0	...	SE	1.79	0	0
1	2010-01-02 01:00:00	148.0	-15	-4.0	...	SE	2.68	0	0
2	2010-01-02 02:00:00	159.0	-11	-5.0	...	SE	3.57	0	0
3	2010-01-02 03:00:00	181.0	-7	-5.0	...	SE	5.36	1	0
4	2010-01-02 04:00:00	138.0	-7	-5.0	...	SE	6.25	2	0
...
43795	2014-12-31 19:00:00	8.0	-23	-2.0	...	NW	231.97	0	0
43796	2014-12-31 20:00:00	10.0	-22	-3.0	...	NW	237.78	0	0
43797	2014-12-31 21:00:00	10.0	-22	-3.0	...	NW	242.70	0	0
43798	2014-12-31 22:00:00	8.0	-22	-4.0	...	NW	246.72	0	0
43799	2014-12-31 23:00:00	12.0	-21	-3.0	...	NW	249.85	0	0

图 5 原始训练集

information of test set:								
	dew	temp	press	wnd_dir	wnd_spd	snow	rain	pollution
0	-16	4	1027	SE	3.58	0	0	128
1	-17	5	1027	SE	7.60	0	0	77
2	-16	4	1027	SE	9.39	0	0	65
3	-16	1	1028	cv	0.89	0	0	79
4	-14	0	1028	NE	1.79	0	0	93
...
341	-23	-2	1034	NW	231.97	0	0	8
342	-22	-3	1034	NW	237.78	0	0	10
343	-22	-3	1034	NW	242.70	0	0	10
344	-22	-4	1034	NW	246.72	0	0	8
345	-21	-3	1034	NW	249.85	0	0	12

图 6 原始测试集

而后，在模型训练之前需要对数据进行预处理，主要步骤如下：

缺失值处理：检查并处理数据中的缺失值，确保数据完整性。

特征编码：将风向（wnd_dir）这一类别变量通过映射方式编码为数值类型，以便于神经网络处理。

NE	SE	NW	cv
0	1	2	3

时间特征提取：将data列转换为datetime格式，并设为索引，以利于时间序列对齐与可视化分析。

数据可视化分析：使用Matplotlib绘制训练集中各特征随时间变化的趋势图；使用Seaborn对测试集各特征进行直方图与核密度分布图分析。

数据归一化：利用MinmaxScaler对污染物浓度及7项气象变量（露点温度、气温、气压、风向、风速、降雪、降雨）进行归一化处理，统一至[0,1]区间，提升训练稳定性。

滑动窗口构建：输入特征为前11小时的7项气象变量，预测目标为当前时间点的PM2.5浓度。构建输入输出对用于LSTM网络训练，实现多变量时间序列监督学习。

[5 rows x 8 columns]								
	pollution	dew	temp	press	wnd_dir	wnd_spd	snow	rain
0	0.128773	0.352941	0.377049	0.654545	0.333333	0.005349	0.0	0.0
1	0.077465	0.338235	0.393443	0.654545	0.333333	0.012219	0.0	0.0
2	0.065392	0.352941	0.377049	0.654545	0.333333	0.015278	0.0	0.0
3	0.079477	0.352941	0.327869	0.672727	1.000000	0.000752	0.0	0.0
4	0.093561	0.382353	0.311475	0.672727	0.000000	0.002290	0.0	0.0

图 7 预处理后训练集

[346 rows x 8 columns]

	pollution	dew	temp	...	wnd_spd	snow	rain
date				...			
2010-01-02 00:00:00	0.129779	0.352941	0.245902	...	0.002290	0.000000	0.0
2010-01-02 01:00:00	0.148893	0.367647	0.245902	...	0.003811	0.000000	0.0
2010-01-02 02:00:00	0.159960	0.426471	0.229508	...	0.005332	0.000000	0.0
2010-01-02 03:00:00	0.182093	0.485294	0.229508	...	0.008391	0.037037	0.0
2010-01-02 04:00:00	0.138833	0.485294	0.229508	...	0.009912	0.074074	0.0

图 8 预处理后测试集

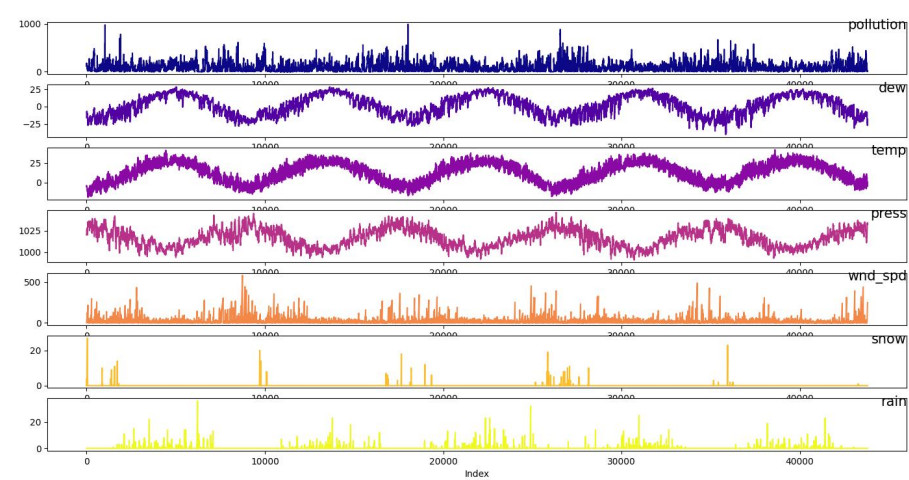


图 9 训练集随时间变化趋势

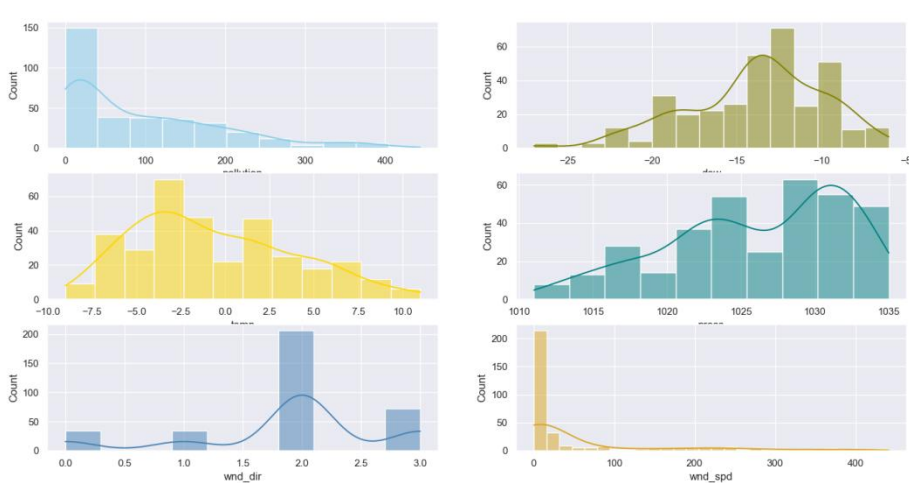


图 10 测试集分布

E2: Training Model Settings

损失函数：采用均方误差（MSE），适合回归问题中度量预测误差。

优化器：使用Adam优化器，学习率设为0.001，具有良好的自适应梯度更新能力。

评估指标：模型训练与测试阶段均采用均方根误差（RMSE）作为性能评估标准，更直观反映预测值与真实值之间的偏差。

训练策略：若验证集损失在10轮内未提升，则提前停止训练以防过拟合，并恢复至表现最好的模型权重；在每轮训练后保存验证损失最优的模型

批量大小：每轮迭代使用32条样本进行权重更新，兼顾训练稳定性与效率。

最大训练轮次：最多训练150轮，在验证集性能不再提升时自动停止。

E3: Results

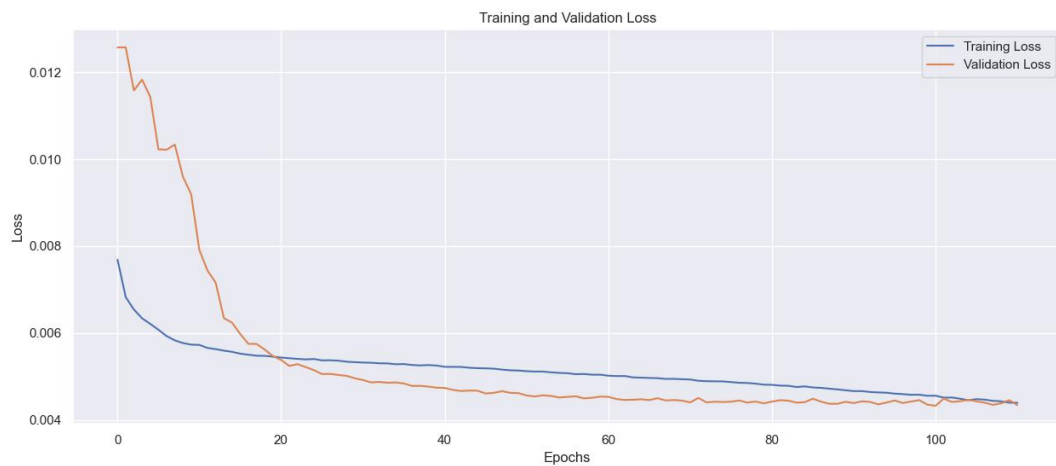


图 11 训练集与验证集损失

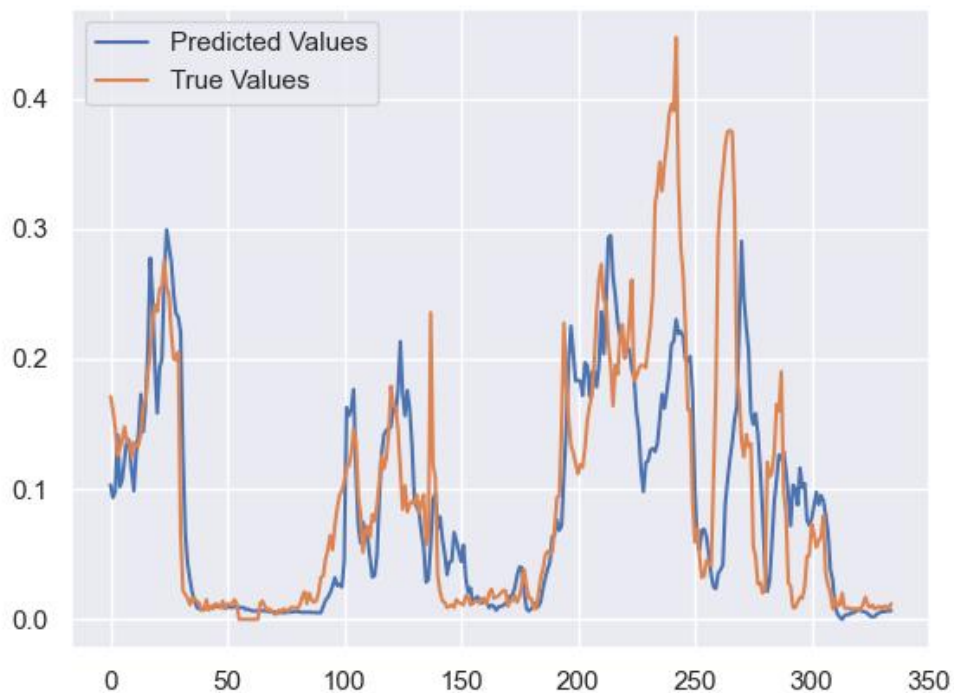


图 12 预测结果与真实结果

Conclusions

本研究基于多变量LSTM网络成功构建了一种用于PM2.5浓度预测的深度学习模型，并在北京地区五年逐小时空气质量与气象数据集上进行了实验验证。通过引入历史时序特征和多种气象因素，模型有效捕捉了污染物浓度随时间演化的非线性变化规律。在合理的数据预处理和特征构造基础上，双层LSTM架构配合Dropout、早停机制以及Adam优化器，不仅增强了模型对时序信息的提取能力，也显著抑制了过拟合风险。实验结果表明，该模型在训练集和验证集上分别达到0.0655和0.0647的均方根误差（RMSE），表现出优良的泛化能力和预测精度，能够较为准确地还原实际PM2.5浓度的变化趋势，尤其在污染波动段依然保持稳定响应。这一结果进一步验证了LSTM在多变量时间序列预测任务中的适用性与优势，为复杂环境下的空气污染建模提供了一条可靠的技术路径，具有重要的现实意义与推广价值。