

doi: 10.3969/j.issn.1007-7375.240247

基于概率分析和数据驱动的高速公路车辆速度预测方法

陈赛飞¹, 傅 惠², 张帅宇³

(1. 广东技术师范大学 电子与信息学院, 广东 广州 510665; 2. 广东工业大学 机电工程学院, 广东 广州 510006;
3. 平高集团有限公司, 河南 平顶山 461670)

摘要: 为改善高速公路车辆速度预测过程中因数据非正态分布而造成的预测误差, 提出一种深度时间卷积网络 (deep temporal convolutional network, DeepTCN) 和 Copula 理论相结合的混合概率性预测方法。为给出考虑多个特征的不确定性预测, 建立了 DeepTCN 框架。根据 DeepTCN 得到的预测结果, 基于预测误差的条件概率分布来拟合适当的 Copula 函数。通过误差补偿得到概率性预测结果。利用在广州市机场高速收集的实际交通数据来验证所提出的混合方法的有效性。数据分析与实验结果表明, 实际数据反映出交通流确实存在高度随机性, 而这种随机性在低车辆密度与高车辆密度时相对较小, 在中等密度时相对较大; 与现有的各种方法相比, DeepTCN 在处理长时间序列信息时存在优越性, 适用于高速公路车辆速度预测场景; 结合 Copula 函数可以在一定程度上补偿数据随机性带来的预测误差, 进一步提高预测精度。

关键词: 深度学习; 数据驱动; 联合概率分布; 交通流预测; Copula 理论

中图分类号: F407.67; U491.1+4

文献标志码: A

文章编号: 1007-7375(xxxx)x-0001-10

A Probability-based and Data-Driven Approach for Highway Vehicle Speed Forecasting

CHEN Saifei¹, FU Hui², ZHANG Shuaiyu³

(1. School of Electronics and Information, Guangdong Polytechnic Normal University, Guangzhou 510665, China;
2. School of Electromechanical Engineering, Guangdong University of Technology, Guangzhou 510006, China;
3. Pinggao Group Co., Ltd., Pingdingshan 461670, China)

Abstract: To cope with the forecast errors caused by non-normal distributions of data in the process of forecasting vehicle speed on highways, this paper proposes a probabilistic forecasting method, which is a hybrid of deep temporal convolutional networks (DeepTCN) and Copula theory. DeepTCN is established first to give a deterministic forecast considering multiple features. Then, an appropriate Copula function is fitted based on the conditional probability distribution of prediction errors according to the obtained forecasts by DeepTCN. Finally, probabilistic results are provided by error compensation. Real traffic data collected on a highway road in Guangzhou, China are utilized to verify the effectiveness of the proposed hybrid method. Data analysis and experimental results show that real-world data reflects significant randomness in traffic flows, with this randomness being relatively mild at both low and high vehicle densities, but more evident at medium densities; DeepTCN demonstrates superior performance in handling long-term time series information compared with various existing methods, making it well-suited for highway vehicle speed forecasting; by incorporating Copula functions, it can compensate for forecast errors caused by data randomness to some extent, further improving forecast accuracy.

Key words: deep learning; data-driven; joint probability distribution; traffic flow forecasting; Copula theory

快速的城市化进程和人口增长加剧了城市交通负担。目前, 不仅是一线大城市, 许多中型城市也面临着与交通相关的问题, 包括交通拥堵加

剧、交通事故频繁、空气质量降低以及化石能源不足等^[1]。针对以上问题, 路网重规划与改建是提高道路通行能力、缓解交通压力的重要途径之一, 但

收稿日期: 2024-06-20

基金项目: 国家自然科学基金资助面上项目 (62273108); 广州市科技计划资助项目 (2023A03J0939)

作者简介: 陈赛飞 (1994—), 女, 浙江省人, 博士, 主要研究方向为系统调度与优化、智能算法、人工智能。

Email: gdutfchen@163.com.

它需要大量的人力、物力、财力和时间资源^[2]。而交通流控制是更具成本效益的方法,在智能交通系统(intelligent transportation systems, ITS)领域有着广泛的应用^[3]。

由于实时交通流控制的决策可信度取决于交通状态的准确性(即模型的输入),作为 ITS 重要的组成部分,短期交通流预测的可靠性能够直接影响交通流控制的效果。在这一背景下,高速公路车辆速度预测作为交通流预测的重要组成部分,发挥着重要的指导作用。得益于检测技术的发展,多类型检测设备,如线圈检测器、电子收费(electronic toll collection, ETC)检测器、车辆以及智能手机的全球定位系统(global positioning system, GPS)等,为高速公路车辆速度预测提供了有力的数据支持^[4],在此基础上,数据驱动的人工智能方法被广泛应用^[5]。相比于传统的基于统计学习的方法^[6],如线性回归、小波分析、混沌理论等,基于机器学习的预测方法具有精度高、泛化能力好的优点。

由于影响高速公路车辆速度的因素较多,如需求时空分布、路网拓扑结构、天气、人的主观因素等,高速公路交通大数据中各特征存在强相关性特点^[4],而这种相关性的确切机制尚未得到充分研究^[7]。在机器学习方法中,深度学习由于其强大的数据挖掘能力,在预测方面显示出巨大潜力。目前,递归神经网络(recurrent neural network, RNN)、卷积神经网络(convolutional neural network, CNN)和图神经网络(graph neural network, GNN)是车辆速度预测中常用的深度学习模型^[7-8]。具体来说,RNN、CNN及其变体,例如长短期记忆(long short-term memory, LSTM)网络、门控递归单元(gated recurrent unit, GRU)、时间卷积网络(temporal convolutional network, TCN),在时间特征提取方面具有良好的性能^[9-10],而基于CNN和GNN的方法分别通过以网格和图形式建模网络来提取空间相关性特征。

现今,自动驾驶等技术的发展对车辆速度预测的准确性和效率性提出了更高的要求,但交通流的高度随机性给这一要求带来巨大挑战。这种随机性主要来源于不确定的需求、异质的路网拓扑结构和多样的相关因素,导致所产生的数据具有非线性、非平稳、多尺度特征的特点^[5]。当数据呈现完全相同时,实际交通状态也可能不同。对此,概率性预

测方法能够较好地应对这一问题。其中,Copula理论作为变量间相关性建模的有力工具,是重要的概率性预测方法之一,在ITS领域已得到广泛应用。Yang等^[11]提出基于多元Copula的车祸计数与碰撞风险建模方法,其中碰撞风险由离散事件计数与连续随机变量来衡量。实验表明,所提出的方法在风险评估与碰撞位置预测方面取得良好的效果。Luan等^[12]定义了行程时间波动率(travel time volatility),利用基于Copula的蒙特卡洛模拟方法刻画路段和路径的行程时间变化,在极端情况下实现了准确的预测结果,并证明该方法在预测路径行程时间可靠性方面的优势。

以上研究促使本文将深度学习技术与Copula理论相结合并应用于高速公路车辆速度预测。由于基于深度学习的交通流预测易受数据输入影响,且深度学习模型存在不可解释性的问题,常产生令人费解的预测误差。Copula理论的引入能够从一定程度上改善因数据非正态分布造成的误差,并通过误差补偿提高预测精度。

本文聚焦于高速公路车辆速度预测,提出一种结合DeepTCN和Copula理论的混合概率性预测方法。首先,建立DeepTCN框架以得到确定性预测结果;在此基础上,选择并拟合适当的Copula函数来补偿DeepTCN产生的预测误差,得到概率性预测结果。利用广州市机场高速的实际数据对所提出的模型进行训练与验证,实验结果验证了所提出方法的有效性。

1 混合预测框架

本文的主要目的是利用相应的随机误差来补偿确定性预测结果,从而得到考虑交通流参数随机特性的概率性预测结果。所提出的预测方法框架如图1所示。

本文所提出的框架包含两个主要模块,分别为DeepTCN模块以及Copula模块,各模块中的具体流程由箭头表示。整体数据集被拆分为4个部分,其中数据集1、2分别用于训练与验证DeepTCN模型;数据集3用于得到最优Copula函数;数据集4用于得到最终的确定性预测结果与给定置信度条件下的概率性预测结果。

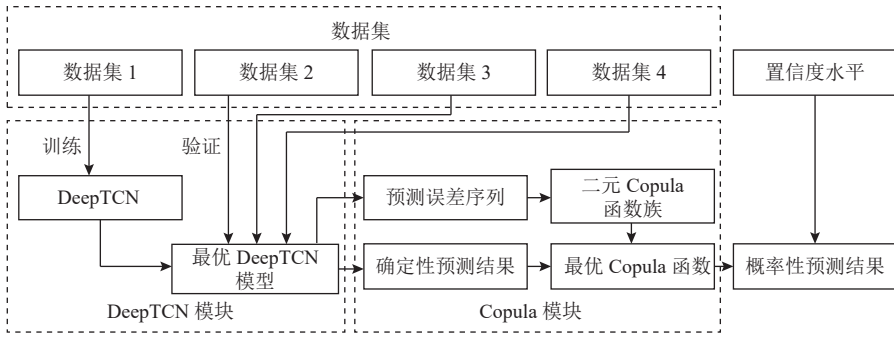


图 1 混合概率性预测框架

Figure 1 Framework of the mixed probabilistic forecasting method

2 基于 DeepTCN 的确定性预测

本研究采用的 DeepTCN 结构如图 2 所示, 主

要由 4 个模块组成, 包括输入模块、编码器、解码器以及输出模块^[13]。在本节中, 将按照这 4 个模块的顺序详细阐述所提出确定性预测的步骤。

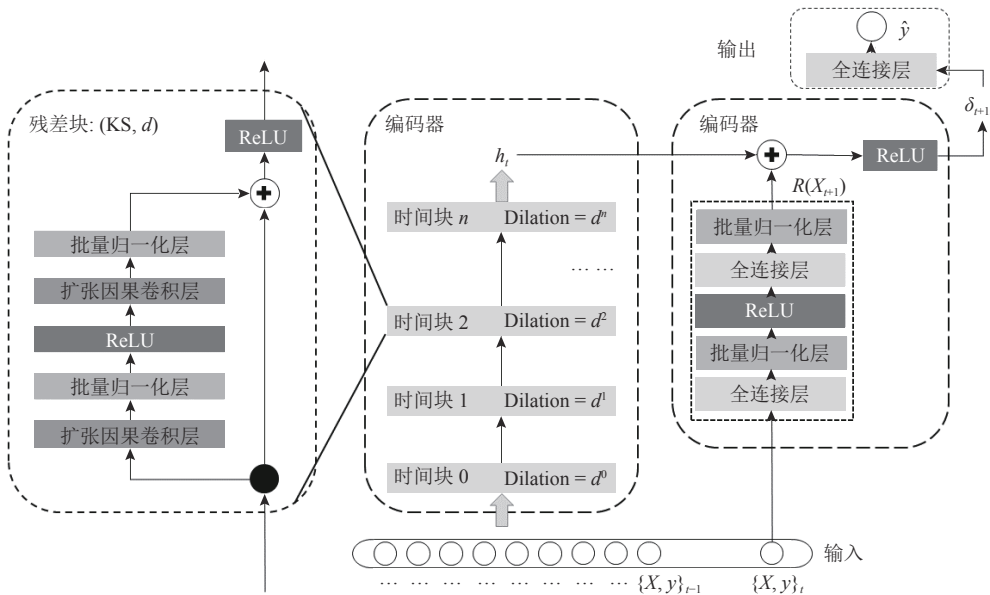


图 2 DeepTCN 结构

Figure 2 Structure of the adopted DeepTCN

在输入模块中, 输入数据集以 $P \times Q$ 矩阵的形式给出, 其中, P 是历史观测长度; Q 代表影响因素的数量。对于每个时间序列 t , 均可以得到一个从 $t-P+1$ 序列开始到 t 序列结束的 $P \times Q$ 矩阵。通过堆叠这些矩阵, 可将二维 $P \times Q$ 矩阵转换为三维矩阵。由于样本数据量巨大, 采用小批量梯度下降算法 (mini-batch gradient descent algorithm) 进行训练, 以平衡模型精度和计算效率。为此, 所得到的三维矩阵需进一步划分为几个小批量, 作为第 2 模块, 即编码器的输入。

在编码器模块中, 所输入的数据信息特征通过扩张因果卷积 (dilated causal convolutions) 进行编

码。通过堆叠多个时间块 (time block), 即残差块 (residual block), 可以得到相当大的感受野, 同时避免数据信息“泄露”于未来时间。基于此, 可调整卷积核尺寸、扩张因子、卷积层数, 来实现对一定时间序列范围内的数据信息的整体感知。式 (1) 给出了输入 x 与卷积核 w 之间在时间序列为 t 时的扩张因果卷积。

$$s(t) = (x * w)(t) = \sum_{k=0}^{KS-1} w(k)x(t-dk) \quad (1)$$

其中, $s(t)$ 为对应的输出; $*$ 代表卷积计算; KS 为卷积核尺寸; d 为扩张因子。

编码器中的每一时间块又包含两个扩张因果卷积层、两个批量归一化层 (batch normalization) 以及一个 ReLU 激活函数层, 其具体排列顺序见图 2 左侧部分。以上网络层组成一个时间块的基本结构, 而第 2 个批量归一化层得到的输出即为该时间块的输入。此后, 再次应用一个 ReLU 层, 以在避免梯度消失的前提下提高模型的鲁棒性和泛化能力。

给定一个 $N \times P \times Q$ 的小批量输入数据矩阵 (即 N 个 $P \times Q$ 的二维矩阵), 通过上述编码器, 即可得到一个 $Q \times 1 \times N$ 的矩阵, 由 h_t 表示。

第 3 模块解码器的输入包含两个部分: 一为编码器的输出 h_t ; 二为外部给定变量 X_{t+1} 。前者代表历史信息, 后者则为一系列能够影响未来交通需求的因子, 如天气、时间序列等, 又称未来协变量。本文中的解码器通过式 (2) 提取 h_t 和 X_{t+1} 的信息。

$$\delta_{t+1} = R(X_{t+1}) + h_t. \quad (2)$$

其中, δ_{t+1} 为解码器输出。

值得注意的是, 式 (2) 中的 $R(\cdot)$ 是一个非线性残差函数, 它将未来协变量 X_{t+1} 映射为预测值 δ_{t+1} 与观测值 h_t 之间的残差。这一过程由残差神经网络变体实现, 如图 2 右侧部分所示。在该变体后应用一个额外的 ReLU 层, 即可得到解码器的输出 δ_{t+1} 。

最后, 通过全连接层的线性降维, 可以得到最终的确定性预测结果 \hat{y} , 这一部分即为所建立 DeepTCN 的输出模块。

本文研究的车辆速度预测属于回归问题, 故采用应用广泛的均方误差损失作为损失函数, 如式 (3) 所示。

$$J = \sum_{j=1}^N (\hat{y}_j - y_j)^2 + L. \quad (3)$$

其中, N 为数据规模, 即样本量; j 为样本索引; \hat{y}_j 与 y_j 分别为预测值与实际值; L 为正则项。本文采用 L1 正则化算法, 即 Lasso 回归, 用于自动选择最重要的特征, 从而降低模型复杂度, 防止过拟合。

3 基于 Copula 理论的预测误差补偿

前文已分析交通流的随机性会反映到交通数据, 进而影响以数据驱动的深度预测模型效果, 那么, 有必要建立预测结果的误差模型以对预

测误差进行统计分析。为此, 考虑利用 Copula 理论构建预测误差的概率分布。

3.1 Copula 函数

Copula 理论由 Sklar 在 1959 年提出^[4]。他指出存在一种名为 Copula 函数的关联函数, 能够描述多个随机变量 (random variable) 的联合累积分布函数 (cumulative distribution function) 与其各自的边缘累积分布函数之间的关系。具体定义如下。

定义 1 给定 n 个随机变量, 它们的联合累积分布函数记为 $F(x_1, \dots, x_n)$, 它们各自的边缘分布函数记为 $F_i(x_i)$, $i = 1, 2, \dots, n$ 。那么, 存在一个如式 (4) 所示的 n 维 Copula 关联函数用于描述 $F(x_1, \dots, x_n)$ 与 $F_i(x_i)$ 的关系。

$$F(x_1, \dots, x_n) = C(F_1(x_1), \dots, F_n(x_n)). \quad (4)$$

其中, $C(\cdot)$ 为 Copula 函数。

定义 2 若边缘分布函数 $F_i(x_i)$, $i = 1, 2, \dots, n$ 为连续的, 那么 Copula 函数 $C(\cdot)$ 唯一确定。

根据定义 1 以及累积分布函数的逆变换, 令 $u_i = F_i(x_i)$, $i = 1, 2, \dots, n$, 则可得到 Copula 函数的表达式为

$$C(u_1, \dots, u_n) = F(F_1^{-1}(u_1), \dots, F_n^{-1}(u_n)). \quad (5)$$

其中, $F_i^{-1}(u_i)$ 为 $F_i(x_i)$ 的反函数。

对于多个连续型随机变量, 基于定义 1、定义 2 以及概率论基本原理, 可将 Copula 函数 $C(\cdot)$ 拆分为各边缘分布的积以及 Copula 密度函数两个部分, 如式 (6) 所示。

$$f(x_1, \dots, x_n) = \frac{\partial^n F(x_1, \dots, x_n)}{\partial x_1 \partial x_2 \dots \partial x_n} =$$

$$\frac{\partial^n C(F_1(x_1), \dots, F_n(x_n))}{\partial x_1 \partial x_2 \dots \partial x_n} = \frac{\partial^n C(u_1, \dots, u_n)}{\partial u_1 \partial u_2 \dots \partial u_n} \cdot \prod_{i=1}^n \frac{\partial F_i(x_i)}{\partial x_i} = c(u_1, \dots, u_n) \cdot \prod_{i=1}^n f_i(x_i). \quad (6)$$

其中, $f(x_1, \dots, x_n)$ 为 $F(x_1, \dots, x_n)$ 对应的概率密度函数 (probability density function); $c(\cdot)$ 为 Copula 密度函数。

3.2 预测误差补偿

本质上来说, 确定性预测结果的误差补偿是为了估计实际值 y 与预测值 \hat{y} 的关系。将二者分别视为连续随机变量 Y 与 \hat{Y} , 根据定义 1 与定义 2, 存在一个唯一的 Copula 函数 $C(\cdot)$ 使之满足式 (7)。

$$H_{Y\hat{Y}}(y, \hat{y}) = C(F_Y(y), G_{\hat{Y}}(\hat{y})). \quad (7)$$

其中, $F_Y(y)$ 和 $G_{\hat{Y}}(\hat{y})$ 分别为随机变量 Y 和 \hat{Y} 的边缘累积分布函数; $H_{Y\hat{Y}}(y, \hat{y})$ 则为随机变量 (Y, \hat{Y}) 的联合累积分布函数。

要利用 Copula 函数来构建联合累积分布函数, 条件之一是各随机变量的边缘分布函数已知。常用方法是对于某随机变量事先假定一种类型的分布, 再利用实际数据拟合得到该分布的具体参数。然而, 在实际预测过程中, 随机变量的分布类型是未知的, 假定的边缘分布类型可能影响最终的联合分布。这也导致本文提出的误差分布拟合无法建立在深度学习框架的损失函数中。为解决这个问题, 本文采用非参数核密度估计方法来拟合适当的边缘分布, 该方法无需事先假定分布类型。具体表达式为

$$\hat{F}_h(x) = \frac{1}{Nh} \sum_{j=1}^N \text{KER} \left(\frac{x - X_j}{h} \right). \quad (8)$$

其中, h 为带宽, 用于控制估计的平滑程度; X_j 为观测的样本数据值; KER(.) 为核函数 (kernel function)。本研究所选择的带宽 h 由考虑估计偏差及其方差的最优平滑规则得到。

由于误差补偿涉及到两个随机变量, 即 Y 与 \hat{Y} , 本文选择五类最常用的二元 Copula 函数来对 Y 与 \hat{Y} 的联合分布进行拟合, 包括两类 Elliptic Copula 函数 (Gaussian Copula、t-Copula) 以及三类 Achimedean Copula 函数 (Gumbel Copula、Clayton Copula、Frank Copula)。每类二元 Copula 函数的参数由最大似然估计法得到。具体过程如下。

根据式 (6), 将 Y 与 \hat{Y} 的联合概率密度函数表示为

$$f_{Y\hat{Y}}(y, \hat{y}) = c(F_Y(y; \theta_y), G_{\hat{Y}}(\hat{y}; \theta_{\hat{y}}), \alpha) \cdot f_Y(y; \theta_y) \cdot g_{\hat{Y}}(\hat{y}; \theta_{\hat{y}}). \quad (9)$$

其中, θ_y 与 $\theta_{\hat{y}}$ 分别为随机变量 Y 与 \hat{Y} 的边缘密度函数涉及的未知参数; $c(F_Y(y; \theta_y), G_{\hat{Y}}(\hat{y}; \theta_{\hat{y}}), \alpha)$ 则为对应的 Copula 密度函数, α 为其未知参数。

通过非参数和密度估计, 根据实际观测到的离散数据序列 (y, \hat{y}) , 可以得到第 j 个样本数据的经验累积分布函数为 $\hat{u}_j = F_Y(y)$ 和 $\hat{v}_j = G_{\hat{Y}}(\hat{y})$ 。在此基础上, Copula 密度函数中的未知参数估计值由式 (10) 得到。

$$\hat{\alpha} = \arg \max \sum_{j=1}^N \text{lnc}(\hat{u}_j, \hat{v}_j, \alpha). \quad (10)$$

估计未知参数值后, 可以基于总体样本数据建立经验 Copula 函数。假设数据规模为 N , 经验 Copula 函数 $\hat{C}(\cdot)$ 由式 (11) 得到。

$$\hat{C}(u_j, v_j) = \frac{1}{N} \sum_{j=1}^N I_{[F_Y(y_j) < \hat{u}_j]} \cdot I_{[G_{\hat{Y}}(\hat{y}_j) < \hat{v}_j]}, \hat{u}_j, \hat{v}_j \in [0, 1]. \quad (11)$$

其中, $I(\cdot)$ 为示性函数, 用于表示事件发生与否。具体来说, 在式 (11) 中, 若事件 $F_Y(y_j) < \hat{u}_j$ 或 $G_{\hat{Y}}(\hat{y}_j) < \hat{v}_j$ 发生, 则 $I_{[F_Y(y_j) < \hat{u}_j]}$ 或 $I_{[G_{\hat{Y}}(\hat{y}_j) < \hat{v}_j]}$ 等于 1, 反之, 等于 0。

在得到经验 Copula 函数之后, 可以计算 Copula 函数 $C(\cdot)$ 与经验 Copula 函数 $\hat{C}(\cdot)$ 各个对应数据点之间的欧氏距离, 如式 (12) 所示。

$$\text{dis} = \sqrt{\sum_{j=1}^N (\hat{C}(u_j, v_j) - C(u_j, v_j))^2}. \quad (12)$$

通过比较所选择的五类二维 Copula 函数的欧氏距离, 可以判断距离最短的 Copula 函数即为最优函数, 并将其用于补偿确定性预测的误差。

建立 Copula 函数后, (Y, \hat{Y}) 的联合分布即可由式 (7) 确定。如前所述, \hat{Y} 的取值由基于 DeepTCN 的确定性预测得到, 即 $\hat{Y} = \hat{y}$ 。在此条件下, 可以得到 Y 的条件概率分布为

$$H_{Y|\hat{Y}}(y|\hat{y}) = \int_{-\infty}^y h_{Y\hat{Y}}(y|\hat{y}) dy = \int_{-\infty}^y \frac{h_{Y\hat{Y}}(y, \hat{y})}{g_{\hat{Y}}(\hat{y})} dy. \quad (13)$$

其中, $H_{Y|\hat{Y}}(y|\hat{y})$ 为 $\hat{Y} = \hat{y}$ 条件下的随机变量 Y 的条件概率分布函数; $h_{Y\hat{Y}}(y, \hat{y})$ 为对应的条件概率密度函数; $h_{Y\hat{Y}}(y, \hat{y})$ 与 $g_{\hat{Y}}(\hat{y})$ 为概率密度函数, 分别与分布函数 $H_{Y\hat{Y}}(y, \hat{y})$ 与 $G_{\hat{Y}}(\hat{y})$ 对应。当 $H_{Y\hat{Y}}(y, \hat{y})$ 与 $G_{\hat{Y}}(\hat{y})$ 确定时, $h_{Y\hat{Y}}(y, \hat{y})$ 与 $g_{\hat{Y}}(\hat{y})$ 即可确定。

假设预测误差 e 为实际值与预测值之差, 即 $e = y - \hat{y}$, 那么误差 e 可由式 (14) 计算得到。

$$e = H_{Y|\hat{Y}}^{-1}(y|\hat{y}) - \hat{y}. \quad (14)$$

综上, 便可对确定性预测进行补偿, 并以概率得到最终的预测结果。

4 实例分析

4.1 数据收集与预处理

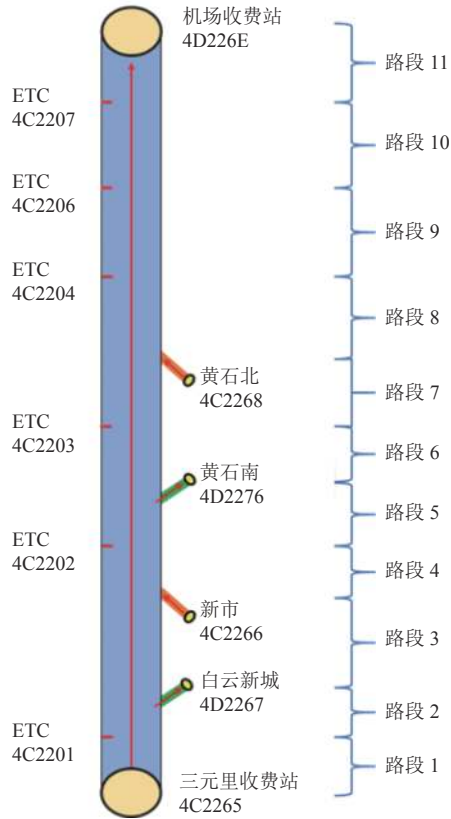
将所提出的方法应用于广州市机场高速车辆平均速度的预测。如图 3 (a) 所示, 广州市机场高速公

路全长 50.47km, 作为交通枢纽, 连接广州市北部地区与新白云国际机场。因此, 该高速路交通压力

较大, 迫切需要一种有效的方法来预测日交通流量, 以便制定相应的控制策略。



(a) OpenStreetMap实际路网



(b) 所选高速路示意图

图 3 实验路网

Figure 3 Test site

本研究为深度学习训练与测试采集了丰富的实际数据, 覆盖时间为 2021 年 9 月 15 日到 2021 年 10 月 15 日。从图 3 可以看出, 数据类型分为两类, 分别为收费站数据 (4C2265, 4D226E, 4D2267, 4D2276, 4C2266, 4C2268) 以及龙门架 ETC 数据 (4C2201, 4C2202, 4C2203, 4C2204, 4C2206, 4C2207)。两类数据包含的字段相同, 分别为检测器 ID、车辆 ID、检测器位置、车辆通过检测器的时间节点。

由于龙门架 ETC 检测器仅能记录装有 ETC 设备的车辆信息, 因此, 收费站数据与龙门架数据中的车辆信息存在差异。由于检测器故障和通信干扰等问题, 龙门架数据存在数据丢失的问题。为避免因数据质量影响预测结果, 首先需要对数据进行预处理。以三元里收费站为起点, 机场收费站为终点这一交通方向为例, 根据检测器位置对车辆数据进行排序, 可以得到具有时间序列的数据集。在此

基础上, 按照以下方法补全缺失的数据。

令 $T(\text{veh}, m)$ 表示车辆 veh 经过第 m 个检测器的时间节点。此外, 令 $D(m, m')$ 为第 m 个检测器与第 m' 个检测器之间的距离。 $T(\text{veh}, m)$ 与 $D(m, m')$ 均可由上述提到的实际数据直接得到。假设车辆 veh 经过第 m_{For} 个检测器与第 m_{Lat} 个检测器的数据记录存在, 且 $m_{\text{For}} < m_{\text{Lat}}$, 而经过第 m 个检测器的数据记录丢失, 那么 $T(\text{veh}, m)$ 可以根据以下 3 种情况估算。

情况 1 若 $m_{\text{For}} < m < m_{\text{Lat}}$, 那么

$$T(\text{veh}, m) = T(\text{veh}, m_{\text{For}}) + (T(\text{veh}, m_{\text{Lat}}) - T(\text{veh}, m_{\text{For}})) \cdot \frac{D(m_{\text{For}}, m)}{D(m_{\text{For}}, m_{\text{Lat}})} \quad (15)$$

情况 2 若 $m < m_{\text{For}} < m_{\text{Lat}}$, 那么

$$T(\text{veh}, m) = T(\text{veh}, m_{\text{For}}) - (T(\text{veh}, m_{\text{Lat}}) - T(\text{veh}, m_{\text{For}})) \cdot \frac{D(m, m_{\text{For}})}{D(m_{\text{For}}, m_{\text{Lat}})} \quad (16)$$

情况 3 若 $m_{For} < m_{Lat} < m$, 那么

$$T(\text{veh}, m) = T(\text{veh}, m_{Lat}) + (T(\text{veh}, m_{Lat}) - T(\text{veh}, m_{For})) \cdot \frac{D(m_{Lat}, m)}{D(m_{For}, m_{Lat})} \quad (17)$$

通过以上轨迹复原处理, 本实验中的 ETC 车辆渗透率由 60% 增加至 90%, 为交通流参数估算提供了数据基础。具体由式 (18) ~ (20) 得到。

$$K(\theta, t) = \frac{TTT(\theta, t)}{l_\theta \cdot \Delta t}; \quad (18)$$

$$Q(\theta, t) = \frac{TTD(\theta, t)}{l_\theta \cdot \Delta t}; \quad (19)$$

$$V(\theta, t) = \frac{TTD(\theta, t)}{TTT(\theta, t)} \quad (20)$$

其中, $K(\theta, t)$ 、 $Q(\theta, t)$ 和 $V(\theta, t)$ 分别为 t 时段内第 θ 个路段的平均车辆密度、流量和速度, 即交通流三参数, 本实验将所选高速公路分为多个路段, 如图 3 (b) 所示; $TTT(\theta, t)$ 和 $TTD(\theta, t)$ 分别表示 t 时段内第 θ 个路段的车辆总体行程时间和总体行驶距离; l_θ 代表第 θ 个路段的长度, 由实际地理位置信息得到; Δt 为数据统计时间步长, 本文设为 5 min。

基于以上统计方法, 本研究共采集 75 178 条数据记录。这些数据被划分为 4 个数据集, 分别为数据集 1、数据集 2、数据集 3、数据集 4(具体用途见图 1), 划分比例为 4 : 1 : 3 : 2。

4.2 交通流随机性分析

选择第 3 个路段 (见图 3 (b)) 对交通流的随机性进行分析。

图 4 为所选路段的交通流基本图 (fundamental diagram), 它反映了平均车辆密度与平均车辆流量之间的关系。以 5min 为统计间隔, 图中每一个红点代表一条数据记录。在此基础之上, 利用三次多项式函数对这一关系进行拟合, 见图中蓝色曲线。拟合得到的曲线走势表明, 在自由流状态下, 路网性能随着车辆密度的增加而提高。然而, 当密度超过一定阈值时, 即 120 veh/km, 路网性能随密度增加而降低, 表明拥堵产生。曲线拐点对应的车辆密度称为关键密度, 它代表最优路网性能所对应的密度, 通常是交通流量控制研究中的关注点之一。值得注意的是, 图 4 中的基本图是高度离散的, 这意味着所选路段中的交通流是高度不均匀的, 证实了交通流存在随机性特点。产生这种不均匀性的原因

可能是车辆时空分布不均、人员驾驶行为不同、检测器分布不均匀等。目前, 主要通过将路网划分为多个区域的方法来消除交通流不均匀带来的影响^[15]。本文对这种随机性特征进行了进一步分析, 以便得到更好的预测结果。

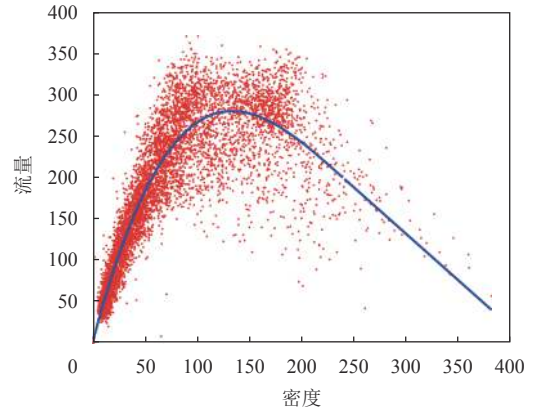
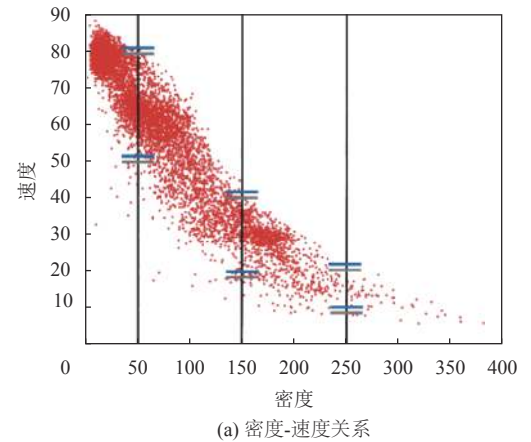


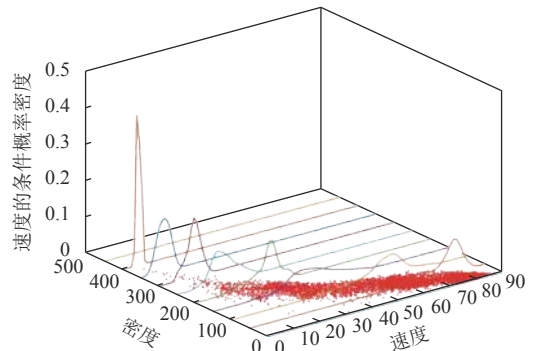
图 4 所选路段的交通流基本图

Figure 4 The fundamental diagram of the selected road section

图 5 展示了所选路段上车辆的平均密度与速度之间的关系。可以明显看出速度与密度之间存在强负相关性。然而, 这一关系的离散程度与红点分布



(a) 密度-速度关系



(b) 速度的条件概率分布

图 5 所选路段的随机基本图

Figure 5 The random fundamental diagram of highway road section 3

的纵向宽度存在非线性相关性。图 5 所呈现的结果表明,低车辆密度(小于 $50 \text{ veh}\cdot\text{km}^{-1}$)与高车辆密度(大于 $250 \text{ veh}\cdot\text{km}^{-1}$)对应的离散程度相对较小,而中等密度(介于 50 到 $250 \text{ veh}\cdot\text{km}^{-1}$ 之间)对应的离散程度较大。一个可能的原因是,在自由流状态下,出行者倾向于以最高限速行驶,各车辆相互作用较小,因此速度差异较小。而在中等密度状态下,由于交通流之间的相互作用变得更大,驾驶行为的多样性被放大,导致离散程度变高(即交通流不均匀性增加)。类似地,当路网处于拥堵状态时,由于相互作用太强,速度只能在很小的范围内变化。

上述现象也可以在图 5 (b) 中得到证实,该子图展现了各给定密度下的速度概率分布。图 5 (b) 显示当密度很小或很大时,速度的概率分布比密度中等时更集中。

通过以上分析,可以证实交通流系统存在时空动态性与随机性,进一步强调了交通流预测的重要性。

4.3 确定性预测结果

确定性预测的数据输入包括由式 (15)~(20) 预处理得到的交通流数据以及对应时间段的天气信息(温度、湿度、降水量)。此外,DeepTCN 的超参数设置为训练迭代次数 $E = 200$,学习率 $\eta = 0.05$,批量大小 $B = 128$,扩张因子 $d = [1,2,4,8,16,32]$ 。

选择 RNN、GRU 和 LSTM 网络与 DeepTCN 比较,根据各方法得到的指标值,来验证本文的确定性预测方法的有效性,所选指标为平均绝对百分比误差(mean absolute percentage error, MAPE)和均方根误差(root mean squared error, RMSE)。表 1 为比较结果。

表 1 各方法预测结果

Table 1 Prediction errors obtained by different methods

方法	MAPE/%	RMSE/($\text{km}\cdot\text{h}^{-1}$)
DeepTCN	5.866	3.923 083
RNN	9.765	5.786 401
LSTM	7.102	4.862 808
GRU	7.273	4.918 785

从表 1 可以看出,本文的 DeepTCN 的预测效果优于其他方法。相比于次优结果(由 LSTM 网络得到),本文提出的方法对于指标 MAPE 提升了 22.4%,对于指标 RMSE 提升了 23.9%。此外,GRU 得到的结果与 LSTM 网络得到的结果相近,而

RNN 得到的结果最差。原因可能在于 RNN 在处理长时间序列信息时存在梯度消失或梯度爆炸的问题。这也进一步证明了 DeepTCN 在处理长时间序列信息时的优越性。

图 6 展示了由 DeepTCN 得到的确定性预测结果,横坐标代表时间,纵坐标代表车辆平均速度。其中红色实线为实际观测值,绿色虚线为预测值。从这两条曲线的走势可以看出预测结果与观测值十分接近,然而部分时间间隔内预测误差较大,如 $t = 90$ 到 $t = 100$,以及 $t = 250$ 到 $t = 264$ 。对此,利用 Copula 理论对得到的预测误差进行补偿,具体结果见 4.4 节。

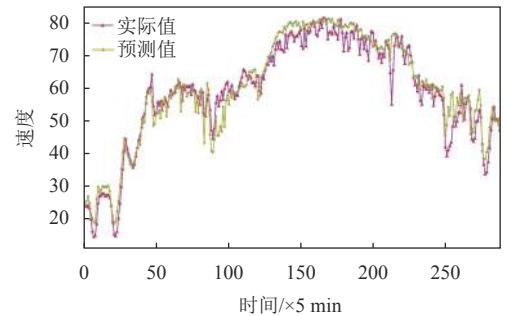


图 6 由 DeepTCN 得到的车辆速度确定性预测结果

Figure 6 Deterministic forecasts of vehicle speed obtained by DeepTCN

4.4 概率性预测结果

在给定概率性预测结果之前,首先需要确定合适的 Copula 函数以描述相关随机变量的联合分布。本文采用车辆速度的预测值与实际值为相关随机变量,以反映预测误差的情况。表 2 对比了三类 Copula 函数的结果,分别为 Gaussian Copula、Gumbel Copula 以及 t-Copula。

表 2 不同 Copula 函数的预测结果对比

Table 2 Forecast comparison between different Copula functions

方法	MAPE/%	RMSEE/($\text{km}\cdot\text{h}^{-1}$)
DeepTCN	0.058 666	3.923 083
DeepTCN + Gaussian Copula	0.051 726	3.637 456
DeepTCN + Gumbel Copula	0.052 523	3.643 571
DeepTCN + t-Copula	0.051 652	3.647 342

从表 2 结果可以看出,DeepTCN 与 Copula 函数的混合方法优于单独的 DeepTCN 方法,验证了 DeepTCN 与 Copula 理论结合的必要性。此外,Gaussian Copula 对于指标 RMSE 最优,而 t-Copula 对于指标 MAPE 最优,在此基础上仍然难以判断二者孰优孰劣。对此,表 3 给出进一步的对比结果。

表 3 不同 Copula 函数的预测效果对比

Table 3 Comparisons of prediction effectiveness using different Copula functions

Copula 函数	欧氏距离	置信度/%	PICP/%	PINAW/%
Gaussian Copula	0.152 7	80	83.96	15.09
		90	90.41	19.41
Gumbel Copula	0.147 0	80	79.33	15.57
		90	87.76	20.88
t-Copula	0.147 5	80	78.34	13.09
		90	88.09	17.18

本文采用 Copula 函数与经验概率分布函数之间的欧氏距离来反映 Copula 函数的拟合效果, 距离越小代表效果越好。从表 3 第 2 列可以看出, 三类 Copula 函数的拟合效果十分相近。此外, 表 3 还给出不同置信度条件下的预测区间覆盖率 (prediction interval coverage probability, PICP) 以及预测区间标准化平均宽度 (prediction interval normalized average width, PINAW)。结果表明, 三类 Copula 函数的指标 PINAW, 但只有 Gaussian Copula 的指标 PICP 超过了给定的置信度水平。这意味着 Gumbel Copula 与 t-Copula 的预测结果可靠性不满足要求。根据以上分析, 本文选择 Gaussian Copula 来进行后续的概率性预测。

图 7 展示了由所提出的混合概率预测方法得到的速度与时间关系图。红色曲线表示所选高速路上车辆平均速度的实际值, 而浅蓝色和深蓝色带分别表示 90% 和 80% 置信度水平条件下的预测区间。结果表明, 两个区间都能很好地覆盖红色曲线, 验证了该方法的有效性。

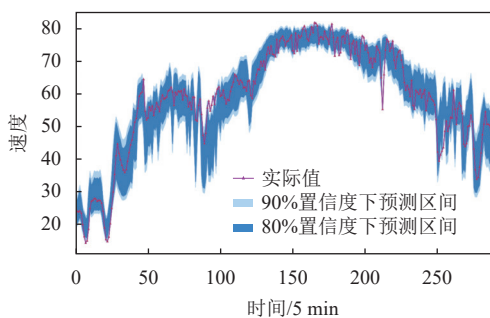


图 7 由本文提出的方法得到的车辆速度概率性预测

Figure 7 Probabilistic forecasts of vehicle speed obtained by the proposed method

5 结论与展望

尽管深度学习技术促进了交通流预测的发展,

但由于交通流具有高度随机性, 在该领域仍然具有很大挑战性。考虑到这些特点, 本文提出一种基于 DeepTCN 和 Copula 理论的交通流混合概率预测方法。首先, 建立 DeepTCN 框架, 并以此模型得到确定性预测结果。然后, 通过合适的 Copula 函数来构建预测值和观测值的联合概率分布。在此基础上, 得到预测误差的条件概率分布。最后, 利用 Copula 学习 DeepTCN 的预测误差规律。结合区间估计理论, 最终得到不同置信度水平下的概率预测结果。基于实际高速公路进行实验, 得出以下结论: 1) 即使在相对简单的高速路路网上, 交通流也存在很高的随机性, 这不仅与时间和空间有关, 还与交通流参数 (如车辆密度和速度) 有关; 2) 与传统的基于深度学习的预测方法相比, DeepTCN 可以提供更高的预测精度; 3) Copula 理论通过考虑不同交通流参数之间的关系, 可以进一步提高预测精度。

本文提出的预测方法目前仅适用于高速公路, 未来将拓展至城市路网。由于城市路网的拓扑结构比高速公路复杂得多, 同时考虑到需求的时空分布, 更应该提取空间相关性。此外, 还将考虑多模式交通 (公交车、私家车、行人、货车等) 对模型的影响, 有必要通过考虑各种交通模式之间的相互作用来进行交通流预测。

参考文献:

- [1] ZHENG F, ZUYLEN VAN H. Urban link travel time estimation based on sparse probe vehicle data[J]. *Transportation Research Part C: Emerging Technologies*, 2013, 31: 145-157.
- [2] OUMAIMA E J, JALEL B O, VERONIQUE V. A stochastic mobility model for traffic forecasting in urban environments[J]. *Journal of Parallel and Distributed Computing*, 2022, 165: 142-155.
- [3] CHEN S, FU H, WU N, et al. Passenger-oriented traffic management integrating perimeter control and regional bus service frequency setting using 3D-pMFD[J]. *Transportation Research Part C: Emerging Technologies*, 2022, 135: 103529.
- [4] JIANG W, LUO J. Graph neural network for traffic forecasting: A survey[J]. *Expert Systems with Applications*, 2022, 207: 117921.
- [5] LIU J, WU N, QIAO Y, et al. Short-term traffic flow forecasting using ensemble approach based on deep belief networks[J]. *IEEE Transactions on Intelligent Transportation Systems*, 2022, 23(1): 404-417.
- [6] 刘甚臻, 马超. 基于小波变换和混合深度学习的短期光伏功率预测[J]. *可再生能源*, 2023, 41(6): 744-749.

- diction based on wavelet transform and hybrid deep learning[J]. *Renewable Energy Resources*, 2023, 41(6): 744-749.
- [7] YU B, LEE Y, SOHN K. Forecasting road traffic speeds by considering area-wide spatio-temporal dependencies based on a graph convolutional neural network (GCN)[J]. *Transportation Research Part C: Emerging Technologies*, 2020, 114: 189-204.
- [8] LV Y, DUAN Y, KANG W, et al. Traffic flow prediction with big data: A deep learning approach[J]. *IEEE Transactions on Intelligent Transportation Systems*, 2015, 16(2): 865-873.
- [9] 王雪琴, 许心越, 伍元凯, 等. 基于混合深度学习模型的城轨短时客流预测[J]. *铁道科学与工程学报*, 2022, 19(12): 3557-3568.
- WANG Xueqin, XU Xinyue, WU Yuankai, et al. Short term passenger flow forecasting of urban rail transit based on hybrid deep learning model[J]. *Journal of Railway Science and Engineering*, 2022, 19(12): 3557-3568.
- [10] WU Y, TAN H, QIN L, et al. A hybrid deep learning based traffic flow prediction method and its understanding[J]. *Transportation Research Part C: Emerging Technologies*, 2018, 90: 166-180.
- [11] YANG D, XIE K, OZBAY K, et al. Copula-based joint modeling of crash count and conflict risk measures with accommodation of mixed count-continuous margins[J]. *Analytic Methods in Accident Research*, 2021, 31: 100162.
- [12] LUAN S, CHEN X, SU Y, et al. Modeling travel time volatility using copula-based Monte Carlo simulation method for probabilistic traffic prediction[J]. *Transportmetrica A - Transport Science*, 2022, 18(1): 54-77.
- [13] CHEN Y, KANG Y, CHEN Y, et al. Probabilistic forecasting with temporal convolutional neural network[J]. *Neurocomputing*, 2020, 399: 491-501.
- [14] SKLAR M. Fonctions de repartition à n dimensions et leurs marges[J]. Paris: Université Paris, 1959, 8: 229-231.
- [15] CHEN S, WU N, FU H, et al. Urban road network partitioning based on bi-modal traffic flows with multiobjective optimization[J]. *IEEE Transactions on Intelligent Transportation Systems*, 2022, 23(11): 20664-20680.

(责任编辑: 郑穗华)