

基于时空特征共享的多源交通数据修复*

王纪禹¹, 陈锐祥¹, 何兆成^{1,2}, 朱依婷^{1,2}, 武智刚¹, 许焕挺¹

1. 中山大学智能工程学院 / 广东省智能交通系统重点实验室, 广东 深圳 518107
2. 鹏城实验室, 广东 深圳 518055

摘要: 本文提出了一种基于多维度特征共享的多层稀疏张量分解模型 (MFS-MSTD)。该模型在 CP (CANDECOMP/PARAFAC) 分解的基础上, 对时空因子矩阵施加低秩正则化; 并采用共享低秩因子矩阵的机制表达多源交通数据之间的互补性, 在非随机缺失或路段级缺失的场景下完成了因子矩阵的梯度更新。实例验证表明: 在速度数据非随机缺失场景下, MFS-MSTD 在 RMSE、MAE 和 MAPE 三个误差指标上相较于基线方法平均降低 17%、21% 和 18%; 在流量数据非随机缺失场景下, RMSE、MAE 和 MAPE 平均降低 52%、54% 和 33%。面对更复杂的路段级缺失场景, MFS-MSTD 的修复性能优于 TGMC-S 和 MTNTF 两个基线模型, 能很好地拟合出未观测路段的交通流量变化。

关键词: 交通数据修复; 多源数据; 共享因子矩阵; CP 分解

中图分类号: U491 **文献标志码:** A **文章编号:** 2097-0137 (2024) 05-0167-10

Joint imputation of multi-source traffic data based on shared multi-dimensional spatiotemporal feature

WANG Jiyu¹, CHEN Ruixiang¹, HE Zhaocheng^{1,2}, ZHU Yiting^{1,2}, WU Zhigang¹, XU Huanting¹

1. School of Intelligent Systems Engineering, Sun Yat-sen University / Guangdong Provincial Key Laboratory of Intelligent Transportation System, Shenzhen 518107, China
2. Pengcheng Laboratory, Shenzhen 518055, China

Abstract: This paper proposes a multi-layer sparse tensor decomposition model based on multidimensional feature sharing (MFS-MSTD). On the basis of CP (CANDECOMP/PARAFAC) decomposition, this model applies low rank regularization to the spatiotemporal factor matrix. The mechanism of sharing a low rank factor matrix is adopted to express the complementarity between multi-source traffic data, and the gradient update of the factor matrix can be completed in non-random missing or segment level missing scenarios. A real experment results show that in the scenario of non-random missing speed data, MFS-MSTD reduces the RMSE, MAE, and MAPE by an average of 17%, 21%, and 18% compared to the baseline method; in the scenario of non-random missing traffic data, RMSE, MAE, and MAPE decreased by an average of 52%, 54%, and 33%. In the face of more complex road segment missing scenarios, the imputation performance of MFS-MSTD is superior to the baseline models TGMC-S and MTNTF, and it can well fit the trend of traffic volume changes in unobserved road segments.

Key words: traffic data imputation; multi-source traffic data; sharing factor matrix; CP decomposition

* 收稿日期: 2024-03-31 录用日期: 2024-04-22 网络首发日期: 2024-07-22

基金项目: 国家重点研发计划(2023YFB4301900); 国家自然科学基金(U21B2090);
中国博士后科学基金(2023M744002)

作者简介: 王纪禹(1998年生), 男; 研究方向: 智能交通系统; E-mail: wangjy365@mail2.sysu.edu.cn

通信作者: 何兆成(1977年生), 男; 研究方向: 智能交通系统; E-mail: hezhch@mail.sysu.edu.cn

全文阅读



ZR20240091

通讯中断以及检测设备故障等因素的存在导致交通检测数据极易发生丢失, 严重影响到智能交通系统中下游任务的准确度 (Said et al., 2022)。早期的研究利用观测值的统计特征来估计发生缺失的交通数据 (Meier et al., 2001; Lint et al., 2005; Ni et al., 2005)。由于相邻路段的交通状态具有明显的相关性, 一些机器学习模型也被引入到交通缺失数据修复任务中。例如: Ma et al. (2020) 提出的 KNN 模型通过度量缺失数据特征与观测数据特征的距离, 然后将距离最近的 K 个数据的均值填补到缺失数据位置。Lu et al. (2018) 提出了一种基于极限学习机自动编码器的缺失数据插补方法。然而, 这些交通数据修复方法的结构较为简单、捕捉复杂和动态交通时空模式的能力较差, 在大规模复杂场景中的修复精度较低。

考虑到交通数据天然具有低秩性的特点, 张量分解被广泛用于高维交通数据的建模与分析 (高远, 2022; Xu et al., 2023)。柏跃龙等 (2019) 使用 Tucker 分解完成对交通流量数据的缺失重构。Wu et al. (2019) 提出了一种改进的 CP 张量分解框架, 该框架结合了多个范数以约束交通数据的平滑性。Chen et al. (2018) 采用 Tucker 分解来描述交通数据中的潜在特征, 并使用截断 SVD 来捕获每个维度中的主要时空特征信息。为避免张量分解受到不确定性因素的干扰, 一些研究对因子矩阵采用正则化的约束以避免模型过拟合。Zhang et al. (2019) 将 F-范数正则化应用于张量分解模型中的因子矩阵, 以提高模型的鲁棒性。Yu et al. (2021) 提出了一种基于 T-Product 的张量分解模型。尽管张量分解模型可以有效捕捉交通数据低秩特性, 但每个时空维度因子矩阵中存在的自相似性却无法得到表达。此外, 现有模型局限于单源交通数据的修复, 在利用多源交通数据的互补优势全面挖掘交通时空模式的问题上, 仍缺少相关研究。

因此, 本文提出了一种多维度特征共享的多层稀疏张量分解方法 (MFS-MSTD, multidimensional feature sharing-multilayer sparse tensor decomposition), 用于同时修复稀疏的交通速度和流量数据。首先, 提出了一种多层稀疏张量分解模型, 该模型基于 CP 分解对每个因子矩阵施加低秩正则约束, 以表征每个因子矩阵的内部自相似性。然后, 基于 MSTD 模型, 采用共享低

秩因子矩阵的机制建立速度和流量数据之间的相关性, 以解决数据大规模缺失而无法计算因子矩阵梯度的问题。最后, 在中国某特大城市快速路的完整速度和流量数据集上进行实验, 验证了所提出方法的有效性和先进性。

1 问题描述

1.1 交通数据的高维组织

本文选择三维张量数据 $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ 为研究对象, 如图 1 所示。其中 I 表示路段数量, J 表示天数, K 表示一天中时间窗数量。张量数据 \mathcal{X} 固定其他维度只保留一个自由度维度时为张量的纤维, 保留两个自由度维度时为张量的切片, 分别记作 $\mathcal{X}(i, j, :)$ 和 $\mathcal{X}(i, :, :)$ 。本文中涉及的其他数学运算符号及其含义如表 1 所示。

表 1 数学符号与含义
Table 1 Symbols and meanings

符号	含义
$\mathcal{X}_{(i)}$	张量 \mathcal{X} 沿第 i 个维度展开得到的矩阵
$\ \mathbf{X}\ _*$	矩阵 \mathbf{X} 的核范数
$\ \mathbf{X}\ _F$	矩阵 \mathbf{X} 的 F-范数
$\langle \mathbf{X}, \mathbf{Y} \rangle$	矩阵内积运算
\odot	Khatri-Rao 积
$\mathbf{I}_{R \times R}$	单位矩阵

1.2 交通数据缺失场景分类

本文主要关注非随机缺失和路段级缺失这两种复杂缺失场景下的交通缺失数据修复任务。

1) 非随机缺失场景: 非随机缺失因其缺失元素呈纤维状分布也被称为纤维级缺失, 如图 1(b) 所示。根据非随机缺失产生原因的不同对应着不同维度的纤维缺失, 例如: $\mathcal{X}(i, j, :) = 0$ 表示在时间维度上的非随机缺失, $\mathcal{X}(:, j, k) = 0$ 则表示在路段维度上的非随机缺失。

2) 路段级缺失场景: 路段级缺失场景是指沿着路段维度失去整个切片数据, 当第 i 条路段发生路段级缺失时, 记作 $\mathcal{X}(i, :, :) = 0$ 。路段级缺失场景是所有交通数据缺失场景中最复杂的一类, 因其缺失了该路段所有的监测信息, 从而很难学习出该路段的高维表征。

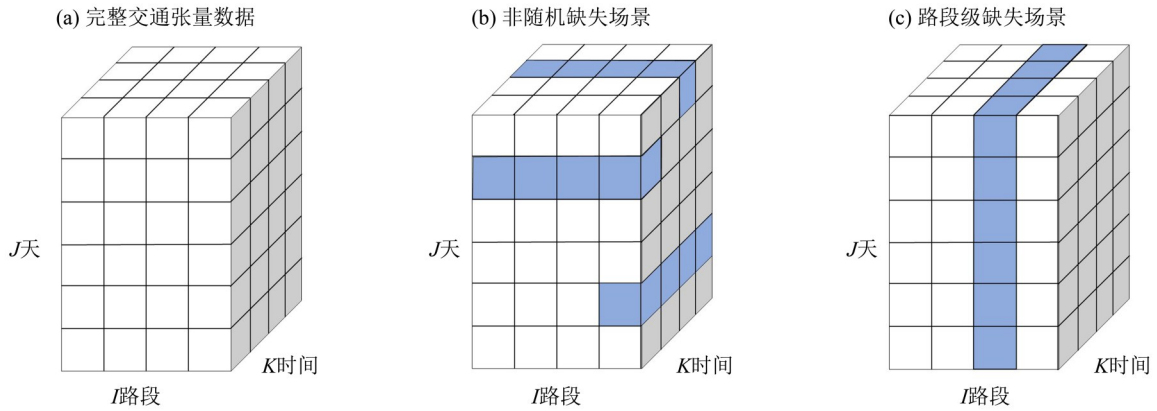


图 1 交通数据缺失场景分类

Fig. 1 Classification of missing traffic data scenarios

2 研究方法

2.1 CP分解

CP分解是张量分解模型中一种经典的高维数据分解结构。对于一个三维交通张量数据 $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ 来说, 当张量秩为 R 时, CP分解可将张量 \mathcal{X} 近似分解为三个时空因子矩阵。

$$\min_{U_1, U_2, U_3} \|\mathcal{X} - [U_1, U_2, U_3]\|_F^2, \quad (1)$$

其中 $U_1 \in \mathbb{R}^{I \times R}$ 表示路段因子矩阵, $U_2 \in \mathbb{R}^{J \times R}$ 表示天因子矩阵, $U_3 \in \mathbb{R}^{K \times R}$ 表示时间因子矩阵。符号“ $[]$ ”为多线性组合运算, 它将三个时空因子矩阵映射为一个完整的交通张量数据。在交通数据修复任务中, CP分解利用已观测的交通数据挖掘出

隐藏的时空特征, 即路段因子矩阵、天因子矩阵和时间因子矩阵。将三个因子矩阵经过多重线性重构后去估计缺失元素。

然而, CP分解只关注交通数据低秩性的特征, 而不考虑路段因子矩阵、天因子矩阵和时间因子矩阵等高维表征的自相似性。因此, 有必要对时空因子矩阵施加低秩正则化的约束, 以全面刻画交通张量数据内在的时空模式。

2.2 多层稀疏张量分解模型

本文提出一个可以同时表达交通数据低秩性和不同维度因子矩阵自相似性的MSTD模型。以一个三维交通张量数据为例, MSTD模型的整体框架如图2所示。

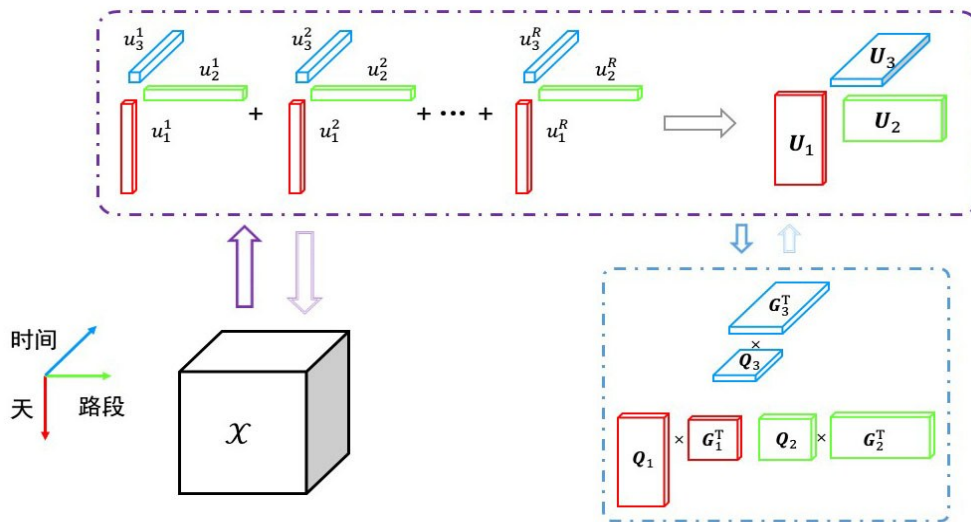


图 2 MSTD模型结构

Fig. 2 The framework of MSTD model

MSTD模型首先计算路段、天和时空间维度的因子矩阵, 利用观测元素挖掘整体交通数据的低秩结构。然后, 继续对每个维度因子矩阵进行低秩分解, 以进一步表达时空因子矩阵的自相似性。MSTD模型对应的优化问题为

$$\min_{U_i, U_2, U_3} \left\| \mathcal{X} - [U_1, U_2, U_3] \right\|_F^2 + \sum_{i=1}^3 \frac{\gamma}{2} \left(\|Q_i\|_F^2 + \|G_i\|_F^2 \right),$$

$$\text{s.t.} \quad U_i = Q_i G_i^T, \quad i = 1, 2, 3,$$

其中 $Q_i \in \mathbb{R}^{I \times R_i}$ 表示第 i 个维度因子矩阵的低秩表征。 $G_i \in \mathbb{R}^{R \times R_i}$ 为第 i 个维度中起缩放作用的矩阵, 它将低秩特征 Q_i 从 R_i 维空间缩放回 R 维空间。参数 γ 为权重因子, 用来平衡因子矩阵低秩项在交通数据修复任务中的重要性。

2.3 多源交通数据联合修复框架

尽管多源交通数据中每种交通数据的统计分布特征不同, 但它们都来自于同一个交通系统并共享相同的时空模式。因此, 本文在MSTD的基础上, 进一步采用共享低秩因子矩阵的机制来表达交通速度和流量数据之间的相互关系。记交通速度张量和流量张量分别为 $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ 和 $\mathcal{Z} \in \mathbb{R}^{I \times J \times K}$, MFS-MSTD模型为

$$\begin{aligned} \min_{U_i, V_i, Q_i, S_i, G_i} & \delta_x \left\| \mathcal{X} - [U_1, U_2, U_3] \right\|_F^2 \\ & + \delta_z \left\| \mathcal{Z} - [V_1, V_2, V_3] \right\|_F^2 + \sum_{i=1}^3 \frac{\gamma_x}{2} \left(\|Q_i\|_F^2 + \|S_i\|_F^2 \right) \\ & + \sum_{i=1}^3 \frac{\gamma_z}{2} \left(\|Q_i\|_F^2 + \|G_i\|_F^2 \right), \end{aligned} \quad (2)$$

$$\text{s.t.} \quad U_i = Q_i S_i^T, \quad V_i = Q_i G_i^T, \quad i = 1, 2, 3,$$

其中 $V_1 \in \mathbb{R}^{I \times R}$ 、 $V_2 \in \mathbb{R}^{J \times R}$ 、 $V_3 \in \mathbb{R}^{K \times R}$ 是流量张量的时空因子矩阵, $U_1 \in \mathbb{R}^{I \times R}$ 、 $U_2 \in \mathbb{R}^{J \times R}$ 、 $U_3 \in \mathbb{R}^{K \times R}$ 是速度张量的时空因子矩阵。 $Q_i \in \mathbb{R}^{I \times R_i}$ 表示第 i 个维度的共享低秩因子矩阵, 也是该维度的低秩表征。 $S_i \in \mathbb{R}^{R \times R_i}$ 和 $G_i \in \mathbb{R}^{R \times R_i}$ 分别为速度张量与流量张量中的非共享因子矩阵, 非共享因子矩阵将共享低秩因子矩阵缩放至每种数据各自的时空因子矩阵。 δ_x 、 δ_z 用于平衡速度张量和流量张量拟合误差的重要性。 γ_x 、 γ_z 用于平衡目标函数中因子矩阵低秩项的重要性。在约束条件中, 不同来源的交通数据通过共享因子矩阵的低秩表征来实现信息互补。

2.4 求解算法

本节将问题(2)拆解为多个子问题, 采用交替

方向乘法求解MFS-MSTD模型中的矩阵变量。问题(2)对应的增广拉格朗日函数为

$$\begin{aligned} \mathcal{L}(U_i, V_i, Q_i, S_i, G_i, W_i, Y_i) & = \delta_x \left\| \mathcal{X} - [U_1, U_2, U_3] \right\|_F^2 + \delta_z \left\| \mathcal{Z} - [V_1, V_2, V_3] \right\|_F^2 \\ & + \sum_{i=1}^3 \frac{\gamma_x}{2} \left(\|Q_i\|_F^2 + \|S_i\|_F^2 \right) + \sum_{i=1}^3 \frac{\gamma_z}{2} \left(\|Q_i\|_F^2 + \|G_i\|_F^2 \right) \\ & + \sum_{i=1}^3 \left(\frac{\sigma_x}{2} \|U_i - Q_i S_i^T\|_F^2 + \langle U_i - Q_i S_i^T, Y_i \rangle \right) \\ & + \sum_{i=1}^3 \left(\frac{\sigma_z}{2} \|V_i - Q_i G_i^T\|_F^2 + \langle V_i - Q_i G_i^T, W_i \rangle \right), \end{aligned} \quad (3)$$

其中 $Y_i \in \mathbb{R}^{I \times R}$ 、 $W_i \in \mathbb{R}^{I \times R}$ 为拉格朗日乘子, σ_x 、 σ_z 为速度和流量等式约束项的惩罚参数。

2.4.1 计算速度张量时空因子矩阵 固定式(3)中的其余变量, 最小化问题(1)可表示为关于速度张量因子矩阵 U_i 的优化问题:

$$\begin{aligned} U_i = \arg \min_{U_i} & \delta_x \left\| \mathcal{X}_{(i)} - U_i (U_{i-1} \odot U_{i+1})^T \right\|_F^2 \\ & + \frac{\sigma_x}{2} \left\| U_i - Q_i S_i^T + \frac{Y_i}{\sigma_x} \right\|_F^2, \end{aligned} \quad (4)$$

其中 $\mathcal{X}_{(i)}$ 为张量 \mathcal{X} 沿第 i 个维度展开所得到的矩阵。计算式(4)的一阶导数并设置其等于零。

$$\begin{aligned} 0 = -2\delta_x & \left[\mathcal{X}_{(i)} - U_i (U_{i-1} \odot U_{i+1})^T \right] (U_{i-1} \odot U_{i+1}) \\ & + \sigma_x \left(U_i - Q_i S_i^T + \frac{Y_i}{\sigma_x} \right), \end{aligned} \quad (5)$$

速度时空因子矩阵通过求解线性方程完成更新, 式(5)可整理为线性方程(6)。

$$\begin{aligned} U_i & \left[2\delta_x (U_{i-1} \odot U_{i+1})^T (U_{i-1} \odot U_{i+1}) + \sigma_x I_{R \times R} \right] \\ & = 2\delta_x \mathcal{X}_{(i)} (U_{i-1} \odot U_{i+1}) + \sigma_x (Q_i S_i^T - \frac{Y_i}{\sigma_x}). \end{aligned} \quad (6)$$

2.4.2 计算流量张量时空因子矩阵 流量张量时空因子矩阵与速度张量时空因子矩阵具有对称性, 虽然对应的变量不同但更新规则类似。更新流量张量时空因子矩阵 V_i 所对应的优化问题为

$$\begin{aligned} V_i = \arg \min_{V_i} & \delta_z \left\| \mathcal{Z}_{(i)} - V_i (V_{i-1} \odot V_{i+1})^T \right\|_F^2 \\ & + \frac{\sigma_z}{2} \left\| V_i - Q_i G_i^T + \frac{W_i}{\sigma_z} \right\|_F^2, \end{aligned} \quad (7)$$

同样, 令式(7)的一阶导数为0得线性方程(8), 进一步整理并求解方程(9)完成流量张量时空因子矩阵 V_i 的更新, 即

$$0 = -2\delta_z \mathcal{Z}_{(i)}(V_{i-1} \odot V_{i+1}) + 2\delta_z V_i (V_{i-1} \odot V_{i+1})^T (V_{i-1} \odot V_{i+1}) + \sigma_z V_i - \sigma_z (Q_i G_i^T - \frac{W_i}{\sigma_z}), \quad (8)$$

$$V_i [2\delta_z (V_{i-1} \odot V_{i+1})^T (V_{i-1} \odot V_{i+1}) + \sigma_z I_{R \times R}] = 2\delta_z \mathcal{Z}_{(i)}(V_{i-1} \odot V_{i+1}) + \sigma_z (Q_i G_i^T - \frac{W_i}{\sigma_z}). \quad (9)$$

2.4.3 计算共享低秩因子矩阵 固定其余矩阵变量, 更新共享低秩因子矩阵 Q_i 时, 需求解优化问题:

$$Q_i = \arg \min_{Q_i} (\gamma_x + \gamma_z) \|Q_i\|_F^2 + \frac{\sigma_x}{2} \left\| U_i - Q_i S_i^T + \frac{Y_i}{\sigma_x} \right\|_F^2 + \frac{\sigma_z}{2} \left\| V_i - Q_i G_i^T + \frac{W_i}{\sigma_z} \right\|_F^2,$$

计算上式的梯度, 则共享低秩因子矩阵 Q_i 更新公式为

$$Q_i (\sigma_z G_i^T G_i + \sigma_x S_i^T S_i + 2(\gamma_x + \gamma_z) I_{R \times R}) = \sigma_z V_i G_i + W_i G_i + \sigma_x U_i G_i + Y_i G_i.$$

2.4.4 计算速度非共享低秩因子矩阵 速度非共享因子矩阵 S_i 起着缩放作用, 将共享低秩因子矩阵缩放至交通速度数据本身适合的数据尺度上, 得

$$S_i = \arg \min_{S_i} \sum_{i=1}^3 \frac{\gamma_x}{2} (\|Q_i\|_F^2 + \|S_i\|_F^2) + \sum_{i=1}^3 \left(\frac{\sigma_x}{2} \|U_i - Q_i S_i^T\|_F^2 + \langle U_i - Q_i S_i^T, Y_i \rangle \right),$$

计算上式的梯度, 通过求解式(10)可完成速度非共享因子矩阵 S_i 的更新, 有

$$S_i (\sigma_x Q_i^T Q_i + 2\gamma_x I_{R \times R}) = \sigma_x U_i^T + Y_i^T Q_i. \quad (10)$$

2.4.5 计算流量非共享因子矩阵 流量非共享因子矩阵与速度非共享因子矩阵在式(11)中同样具有对称性, 根据 S_i 的更新方式可以直接得到 G_i 的更新方式, 即

$$G_i (\sigma_z Q_i^T Q_i + 2\gamma_z I_{R \times R}) = \sigma_z V_i^T + W_i^T Q_i. \quad (11)$$

2.4.6 更新拉格朗日乘子 对拉格朗日乘子 $\{W_i, Y_i\}_{i=1, 2, 3}$ 按式(12)进行更新:

$$\begin{cases} W_i = W_i + \sigma_z (V_i - Q_i G_i^T), \\ Y_i = Y_i + \sigma_x (U_i - Q_i S_i^T). \end{cases} \quad (12)$$

3 研究结果

3.1 实验设置

3.1.1 数据集 选择中国某特大城市一条由西向东的快速路作为实例, 如图3所示。该快速路全长12.3 km, 根据出入口匝道的布局可以将其划分为18条路段。数据集的采集时间为2019年12月2日到2019年12月29日, 共计28天。采集得到的速度与流量数据以每5 min为一个时间窗进行集计, 一天内共有288个时间窗。



图3 实例验证的场景

Fig. 3 The scene for instance verification

3.1.2 参数设置 实验的硬件环境为 Intel(R) Core(TM) i5-1155G7 2.50 GHz 的 CPU, 所有算法都基于 Python 3.7、Numpy 实现。MFS-MSTD 模型的参数设置为: $\delta_x = 10^{-2}$ 、 $\delta_z = 10^{-5}$; 因子矩阵低秩项的权重参数 $\gamma_x = 10^{-8}$ 、 $\gamma_z = 10^{-8}$; 惩罚因子 σ_x 和 σ_z 的初始值是 10^{-7} , 随着迭代次数的增加每次提高 1.05 倍; 最大迭代次数 K 和最小迭代误差 ϵ 分别设置为 100 和 10^{-5} 。

3.1.3 基线方法与评价指标 考虑到 MFS-MSTD 模型是基于 CP 分解框架而开发的多源交通数据联合修复模型, 因此本文选择了一些广泛应用的 CP 分解结构作为基线方法。在非随机缺失和路段级缺失两种场景中, 比较了 MFS-MSTD 与 KNN、CP

(Acar et al., 2011)、CP_Norm (Zhang et al., 2019)、TRTF (Chen et al., 2022) 和 MSTD 等方法的交通数据修复性能。采用均方根误差 (RMSE)、平均绝对误差 (MAE) 和平均绝对百分比误差 (MAPE) 三个精度评价指标对模型性能进行分析。

3.2 非随机缺失场景下的性能比较

表 2-3 记录了高缺失下的修复性能, 图 4 则绘制了缺失率从 20% 到 80% 时的实验结果。在非随机缺失场景下, 流量张量数据的修复难度要远远高于速度张量数据的修复难度。相同缺失率下, 同一方法在速度任务中的估计误差低于流量任务的估计误差。各方法的修复精度随着缺失率的增加而出现明显下降。

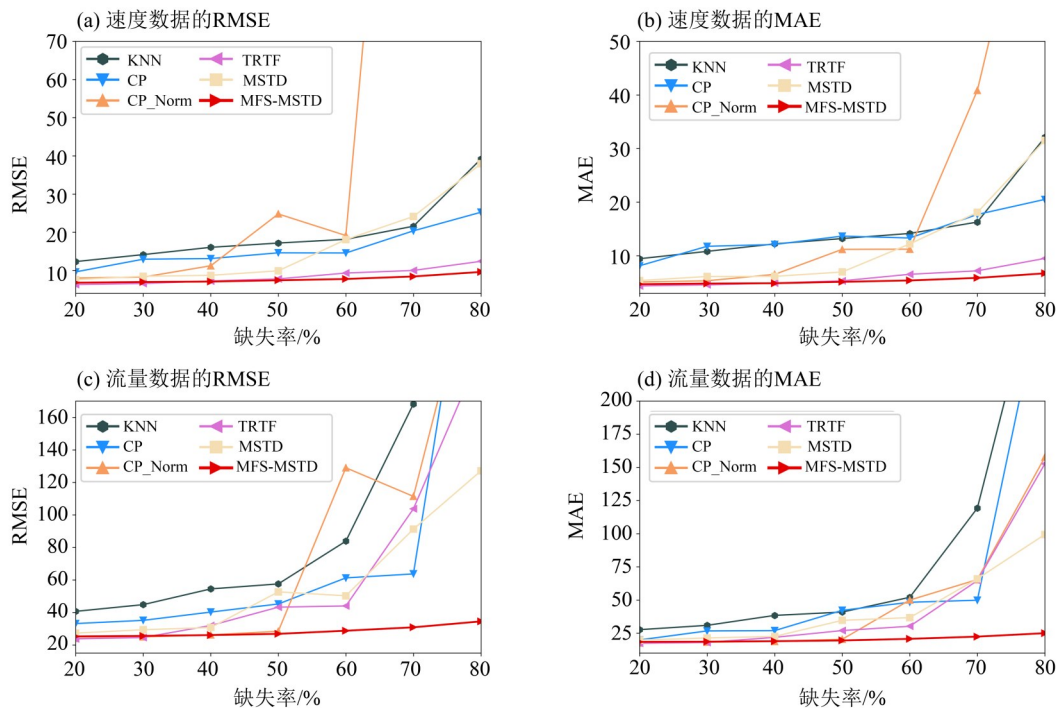


图4 非随机缺失场景下的修复性能

Fig. 4 Imputation performance in non-random missing scenarios

在高缺失率的实验场景中, 本文所提出的 MFS-MSTD 模型的各个指标都达到最佳性能。在交通速度任务的修复实验中, MFS-MSTD 相比于基线方法在 RMSE、MAE 和 MAPE 三个指标上分别降低了 17%、21% 和 18%。而在交通流量任务的修复实验中, MFS-MSTD 的准确度显著提高。不仅在 RMSE、MAE 和 MAPE 三个指标上分别平均降低了 52%、54% 和 33%, 特别在 80% 缺失率的条件下 RMSE、MAE 和 MAPE 降低了 72%、74% 和 45%。在 50% 缺失率下, 对速度数据和流量数

据的修复结果如图 5-6 所示。其中, 淡绿色区域表示发生数据缺失的时间段, 蓝色线为真实值, 红色线代表 MFS-MSTD 方法的估计结果。

与 MSTD 相比, MFS-MSTD 在交通速度和流量的修复任务上都得到了显著提升。这表明, 将速度和流量这两个相互关联的数据集进行联合修复是可行有效的策略。共享低秩因子矩阵机制在多源交通数据的协同修复过程中至关重要, 其能够帮助缺失了大部分元素的维度表征实现梯度更新。在缺失率较低时, KNN 能够得到较为准确的

表2 在非随机缺失场景下各方法在速度张量上的修复性能

Table 2 The imputation performance of various methods on speed tensor in non-random missing scenarios

缺失率	误差指标	KNN	CP	CP_Norm	TRTF	MSTD	MFS-MSTD
60%	RMSE	18.13	14.55	19.13	9.30	18.06	7.75
	MAE	14.16	13.31	11.29	6.54	11.43	5.39
	MAPE	39.61%	32.39%	31.02%	18.98%	29.18%	15.75%
70%	RMSE	21.56	20.35	216.14	9.98	24.10	8.40
	MAE	16.28	17.70	40.82	7.19	18.07	5.87
	MAPE	43.55%	42.55%	121.59%	20.58%	40.76%	17.09%
80%	RMSE	39.16	25.23	259.78	12.37	37.93	9.59
	MAE	32.09	20.51	93.63	9.50	25.33	6.72
	MAPE	72.04%	46.51%	236.42%	25.33%	66.82%	19.65%

表3 在非随机缺失场景下各方法在流量张量上的修复性能

Table 3 The imputation performance of various methods on volume tensor in non-random missing scenarios

缺失率	误差指标	KNN	CP	CP_Norm	TRTF	MSTD	MFS-MSTD
60%	RMSE	83.79	61.20	128.93	43.95	50.20	28.71
	MAE	52.01	48.25	49.91	30.16	36.70	20.68
	MAPE	44.88%	89.81%	89.81%	42.79%	68.77%	37.45%
70%	RMSE	168.16	63.64	111.31	103.67	91.186	30.81
	MAE	119.21	49.80	65.35	64.55	65.947	22.36
	MAPE	193.76%	101.82%	69.67%	66.22%	138.51%	37.57%
80%	RMSE	326.78	324.60	274.01	199.98	127.02	34.45
	MAE	318.90	291.77	157.67	152.66	99.39	24.95
	MAPE	211.25%	99.99%	211.90%	91.36%	193.92%	49.94%

估计,但随着缺失率的增加其性能显著下降。由于只能从观察到的交通数据中获得低秩先验,CP分解在处理非随机缺失场景时性能受到限制。非随机缺失场景对CP_Norm和TRTF模型来说都具有挑战性,这两个模型依赖于因子矩阵稀疏性或时间正则项的约束。从实验结果来看,相比于朴素的CP分解以及CP_Norm,TRTF具有明显优势,这说明时序关系也是交通数据的一个重要规律。MSTD对时空因子矩阵采用的低秩约束相比于因子矩阵稀疏化也更有优势。

3.3 路段级缺失场景下的性能比较

交通流量这一参数的获取依赖于检测设备的实时监测,但受到投资成本的限制,一些路段上并未安装检测器从而导致了交通流量呈现路段级缺失。值得注意的是,相比于流量数据,速度数据的获取来源更加广泛,其可以通过路网上广泛分布的浮动车获取。因此,交通流量参数往往呈

现路段级缺失状态,而交通速度参数则为随机缺失状态。上一节中采用的CP、CP_Norm、TRTF和MSTD等都是针对于单源交通缺失数据的修复方法,无法很好处理路段级缺失的流量数据。本节将检验MFS-MSTD模型在流量数据呈路段级缺失下的修复性能。

本文设置了三种路段级流量数据缺失场景,如表4所示,并将MFS-MSTD模型与TGMC-S(Zhu et al., 2022)和MTNTF(Zhang et al., 2020)模型进行修复精度比较。TGMC-S方法通过交通速度数据计算出每条路段的相似度矩阵,并将其作

表4 三种路段级缺失场景

Table 4 Three scenarios for road segments missing

场景编号	缺失路段的集合	实际缺失率
1	r_6	7.365%
2	r_4, r_{13}	12.748%
3	r_1, r_7, r_{14}	18.271%

为先验信息提供给流量修复任务,通过矩阵补全框架完成流量缺失数据的修复。MTNTF 则是一个多任务学习的深度学习框架,通过学习流量数据和速度数据各自的非线性表征修复数据。在路段级缺失场景中,流量数据通过共享速度数据的路段表征完成缺失值的修复。

表 5 展示了三种方法在路段级缺失实验场景中的表现。在三个实验场景下, MFS-MSTD 模型都取得了优异的修复结果;在场景 2 和场景 3 中的精度指标有明显提升。在场景 2 中, MFS-MSTD 模型的 RMSE 为 30.657、MAE 为 23.746,相较于基线方法误差指标分别降低了 23% 和 27%。在场景 3 中, MFS-MSTD 模型的 RMSE 为 67.905,相较于

基线方法误差指标降低了 10%。在多源交通数据中,由于每种数据的统计分布不一致,并且速度与流量数据呈现的是非线性关系,因此相同路段所采集到的流量和速度数据的相似度也可能不同,导致 TGMC-S 方法在不同场景中的表现差异较大。因其多任务深度学习框架的优势, MTNTF 在这三种路段级缺失场景中的表现都比较鲁棒,但流量数据的统计特征相较于速度有明显差异。MFS-MSTD 模型共享的是低秩因子矩阵,这是交通数据中最本质的时空模式规律。基于此条件的共享机制可以避免受到统计特征带来的偏置,从而让模型更加鲁棒。

表 5 路段级缺失场景下的修复性能

Table 5 Imputation performance in road segment missing scenarios

场景编号	TGMC-S		MTNTF		MFS-MSTD	
	RMSE	MAE	RMSE	MAE	RMSE	MAE
1	33.401	26.345	30.571	23.379	31.406	23.774
2	51.966	36.939	40.156	32.863	30.657	23.746
3	115.975	88.713	76.262	56.676	67.905	54.710

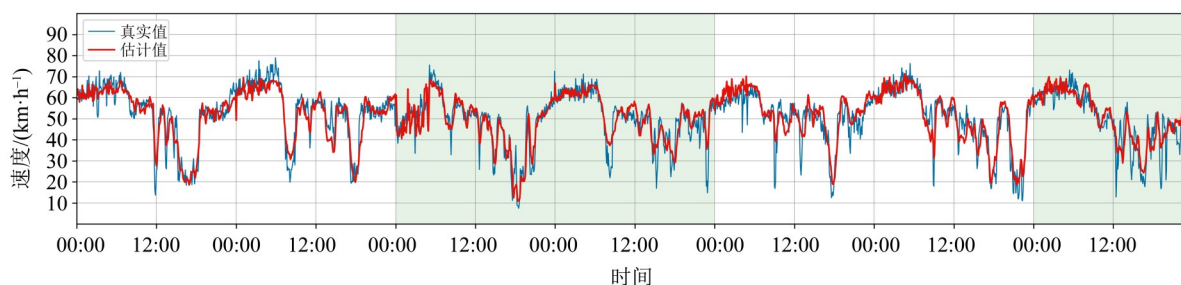


图 5 在速度数据 50% 非随机缺失场景下的修复结果(MFS-MSTD)

Fig. 5 The imputation results in the scenario of non-random missing of 50% speed data(MFS-MSTD)

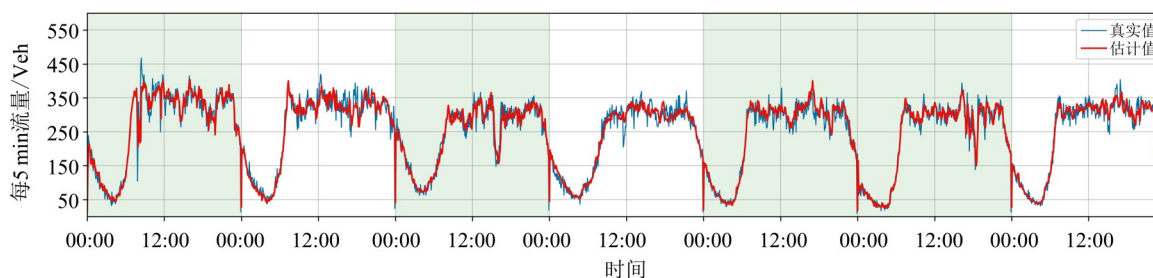


图 6 在流量数据 50% 非随机缺失场景下的修复结果(MFS-MSTD)

Fig. 6 The imputation results in the scenario of non-random missing of 50% volume data(MFS-MSTD)

为了直观展示 MFS-MSTD 模型在路段级缺失场景下的表现,图 7 绘制了场景 1 中 r_6 这一未观测路段一周的修复结果。从图中可以看出,即使在

路段完全没有观测信息的条件下,通过速度数据提供的低秩因子矩阵, MFS-MSTD 模型也能很好地拟合出未观测路段的交通流量动态变化趋势。

3.4 共享低秩因子矩阵的可解释性

共享低秩因子矩阵本质上是时空因子矩阵的左分解矩阵, 在物理意义上其表征时空因子矩阵的低秩自相似特征。由图 8 中路段 r_8 和 r_9 的速度-流量基本图可知, 这两条路段有着极其相似的速度-流量散点分布。因此, 共享路段因子矩阵体现的是路段基本图特征。从路段维度来看, 如图 9(a) 所示, 在共享路段因子矩阵中路段 r_8 和 r_9 呈现明显相似性。从天维度来看, 共享天因子矩阵

明显呈现以 7 天为周期的分布特征, 如图 9(b) 所示。最后一周的相似度比较混乱, 可能是由于周中有圣诞节从而改变了人们的出行模式。因此, 共享天因子矩阵刻画的是以星期为周期的交通出行模式。如图 9(c) 所示, 从时间维度来看, 共享时间因子矩阵具有明显的低秩特点, 在一天内以晚高峰与其他时段的交通模式为主要特征。因此, 对时间维度来说, 共享时间因子矩阵体现的是晚高峰和平峰的交通模式。

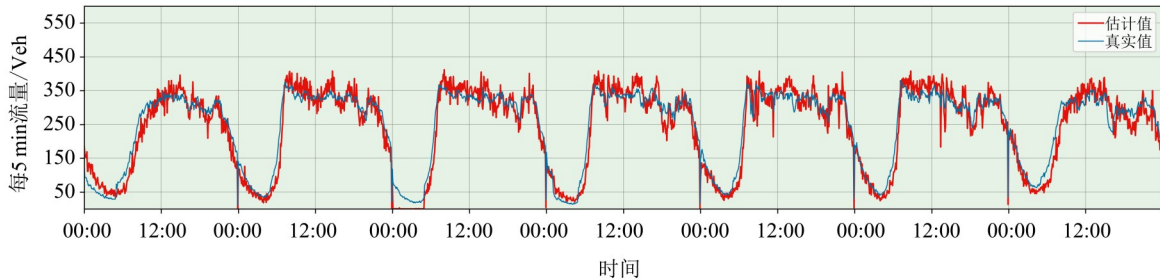


图7 场景1中路段 r_6 的流量数据修复结果

Fig. 7 The imputation results of volume data for section r_6 in scenario 1

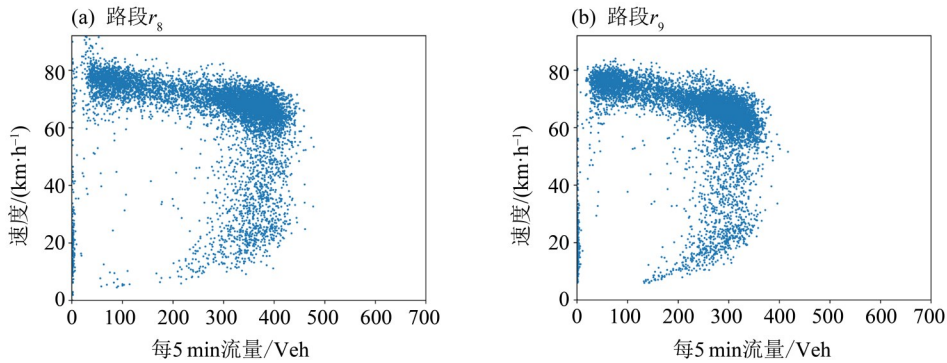


图8 路段速度-流量基本图

Fig. 8 Speed-volume fundamental diagram of road

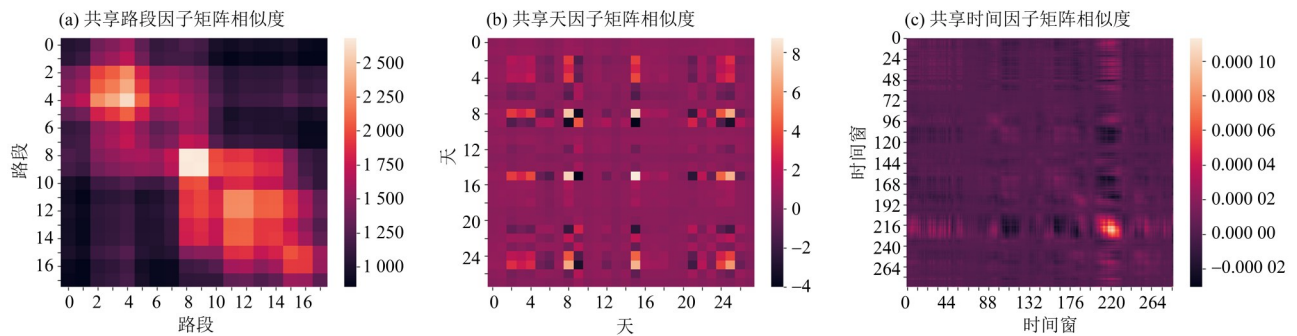


图9 共享低秩因子矩阵相似度

Fig. 9 Shared low-rank factor matrix similarity

4 结 论

本文提出了一种基于共享低秩因子矩阵机制的多维度特征共享的多层稀疏张量分解方法,以同时修复交通速度和流量数据;并以中国某特大城市的一条快速路为例,开展了非随机缺失场景和路段级缺失场景的交通数据修复实验。实验结果表明:在速度数据非随机缺失场景下,

MFS-MSTD 相较于基线方法在 RMSE、MAE 和 MAPE 三个误差指标上平均降低 17%、21% 和 18%。流量数据非随机缺失场景下 RMSE、MAE 和 MAPE 平均降低 52%、54% 和 33%。面对更复杂的路段级缺失场景, MFS-MSTD 取得了优于包括深度学习方法在内的先进模型的精度,并能很好地拟合出未观测路段的交通流量变化趋势。

参考文献:

- 柏跃龙,彭理群,祁钰茜,等,2019.检测数据缺失条件下的交通流估计方法研究[J].交通信息与安全,37(2):99-106.
- 高远,2022.基于稀疏数据的区域路网交通状态识别与预测技术研究[D].北京:北京交通大学.
- ACAR E, DUNLAVY D M, KOLDA T G, et al, 2011. Scalable tensor factorizations for incomplete data [J]. *Chemomet Intell Lab Syst*, 106(1): 41-56.
- CHEN X Y, HE Z C, WANG J W, 2018. Spatial-temporal traffic speed patterns discovery and incomplete data recovery via SVD-combined tensor decomposition [J]. *Transp Res Part C Emerg Technol*, 86: 59-77.
- CHEN Y P, ZHANG Y Q, XIA H, et al, 2022. A hybrid tensor factorization approach for QoS prediction in time-aware mobile edge computing [J]. *Appl Intell*, 7(52): 8056-8072.
- LINT J W C, HOOGENDOORN S P, ZUYLEN H J, 2005. Accurate freeway travel time prediction with state-space neural networks under missing data [J]. *Transport Res C-Emerg Technol*, 13(5/6): 347-369.
- LU C B, MEI Y, 2018. An imputation method for missing data based on an extreme learning machine auto-encoder [J]. *IEEE Access*, 6: 52930-52935.
- MA Z F, TIAN H P, LIU Z C, et al, 2020. A new incomplete pattern belief classification method with multiple estimations based on KNN[J]. *Appl Soft Comput*, 90: 106175.
- MEIER J, WEHLAN H, 2001. Section-wise modeling of traffic flow and its application in traffic state estimation [C]// *IEEE Trans Intell Transp Syst*. IEEE: 440-445.
- NI D H, LEONARD J D, GUIN A, et al, 2005. Multiple imputation scheme for overcoming the missing values and variability issues in ITS data[J]. *J Transp Eng*, 12(131): 931-938.
- SAID A B, ERRADI A, 2022. Spatiotemporal tensor completion for improved urban traffic imputation [J]. *IEEE Trans Intell Transp Syst*, 7(23): 6836-6849.
- WU Y, TAN H, LI Y, et al, 2019. A fused CP factorization method for incomplete tensors [J]. *IEEE Trans Neural Netw Learn Syst*, 30(3): 751-764.
- XU X Q, LIN M W, LUO X, et al, 2023. HRST-LR: A hessian regularization spatio-temporal low rank algorithm for traffic data imputation [J]. *IEEE Trans Intell Transp Syst*, 10(24): 11001-11017.
- YU G H, WANG L Q, WAN S C, et al, 2021. Tensor factorization with total variation for internet traffic data imputation[J]. *Pacific J Optim*, 17(3): 486-505.
- ZHANG H, CHEN P, ZHENG J F, et al, 2019. Missing data detection and imputation for urban ANPR system using an iterative tensor decomposition approach [J]. *Transp Res Part C Emerg Technol*, 107: 337-355.
- ZHANG Z C, LI M, LIN X, et al, 2020. Network-wide traffic flow estimation with insufficient volume detection and crowdsourcing data [J]. *Transp Res Part C Emerg Technol*, 121: 102870.
- ZHU Y T, WANG J Y, WANG J B, et al, 2022. Multitask neural tensor factorization for road traffic speed-volume correlation pattern learning and joint imputation[J]. *IEEE Trans Intell Transp Syst*, 23(12): 24550-24560.

(责任编辑 王海蓉)