# Correlation and causality between traffic congestion and the built environment: a case study in New York city

Weihua Huan, Songnian Li, Xintao Liu, Hangbin Wu, Mi Diao, Hao Li, A. Yair Grinberger, Haobing Liu, Chun Liu & Wei Huang

View supplementary material ⬀

Published online: 07 Sep 2025.

Submit your article to this journal ⬀

Article views: 935

View related articles ⬀

View Crossmark data ⬀

# Correlation and causality between traffic congestion and the built environment: a case study in New York city

Weihua Huan[a,b], Songnian Li[c], Xintao Liu [b], Hangbin Wu[a,d], Mi Diao[e], Hao Li[f,g], A. Yair Grinberger[h], Haobing Liu[d,i], Chun Liu[a,d] and Wei Huang[a,c,d]

[a]College of Surveying and Geo-informatics, Tongji University, Shanghai, People's Republic of China; [b]Department of Land Surveying and Geo-Informatics, The Hong Kong Polytechnic University, Kowloon, Hong Kong SAR, People's Republic of China; [c]Department of Civil Engineering, Toronto Metropolitan University, Toronto, Canada; [d]Urban Mobility Institute, Tongji University, Shanghai, People's Republic of China; [e]College of Architecture and Urban Planning, Tongji University, Shanghai, People's Republic of China; [f]Department of Geography, National University of Singapore, Singapore; [g]Department of Aerospace and Geodesy, Professorship for Big Geospatial Data Management, Technical University of Munich, Munich, Germany; [h]Department of Geography, The Hebrew University of Jerusalem, Jerusalem, Israel; [i]College of Transportation, Tongji University, Shanghai, People's Republic of China

**ABSTRACT**

Traffic congestion is significantly affected by the built environment. Existing studies predominantly examine this through correlation analysis, overlooking causal mechanisms. This omission leads to unreliable feature selection in policy models and hinders evidence-based interventions. To address this, this study proposes a three-stage causal framework that rigorously assesses built environment impacts. The first stage identifies statistically significant correlations using multivariable least squares regression. The second stage applies five causal inference models – Granger causality, structural equation model, causal forest, causal impact, and convergent cross mapping – to uncover causality. The third stage assesses how the identified causal factors shape congestion patterns in perpetually congested roadways (PCRs). Applied to New York City (NYC), the United States, the results reveal 19 correlated and 11 causal impacts. Our key findings include: (1) Transit accessibility is the most robust causal factor, while built environment diversity exhibits time-dependent variability; (2) traffic light design demonstrates bidirectional causality with congestion; (3) PCRs exhibit four distinct spatiotemporal patterns, with bridge-related congestion having the most consistent impact. These results yielded policy recommendations for NYC transportation planning: (i) improve the first-and-last-mile connectivity through micro-mobility; (ii) deploy artificial intelligence-driven adaptive traffic signals; (iii) expand the capacity of critical bridge corridors near PCRs.

## 1. Introduction

Amidst the broader context of sustainability goals articulated in the 2030 Agenda for the United Nations' Sustainable Development Goals (SDGs), traffic congestion has emerged as a pressing scientific concern (Kaiser and Deb 2025; Saberi et al. 2020). Yet, addressing congestion remains a multifaced complex challenge. For instance, the unprecedented pace of urbanization has significantly exacerbated urban congestion (Su et al. 2025). Moreover, congestion is also shaped by the rise of Digital Earth technologies and their use in transportation (Jiang et al. 2022; Li et al. 2019). All these factors underscore the urgency of mitigating congestion – a critical objective for aligning urban development over the coming decade with the SDGs. Achieving this requires uncovering the underlying causal drivers of congestion.

Existing studies have demonstrated that the built environment is a central factor in directly shaping travel behavior (Benito-Moreno, Carpio-Pinedo, and Lamíquiz-Daudén 2025; Gao et al. 2023; Liu et al. 2025; Tracy et al. 2011; Zhang, Sun, and Zegras 2021). In exploring the correlative relationship between the built

environment and behavior, one branch of research classifies physical features using five 'Ds' schema (Cervero and Kockelman 1997; Li et al. 2024; Wang and Zhou 2017): Diversity (the degree of land use mixture within a defined area), Density (concentration of population or physical structures within a given area), Design (Configuration and physical attributes of transportation infrastructure), Distance to transit (proximity to transit hubs), and Destination accessibility (access to transportation nodes or destinations).
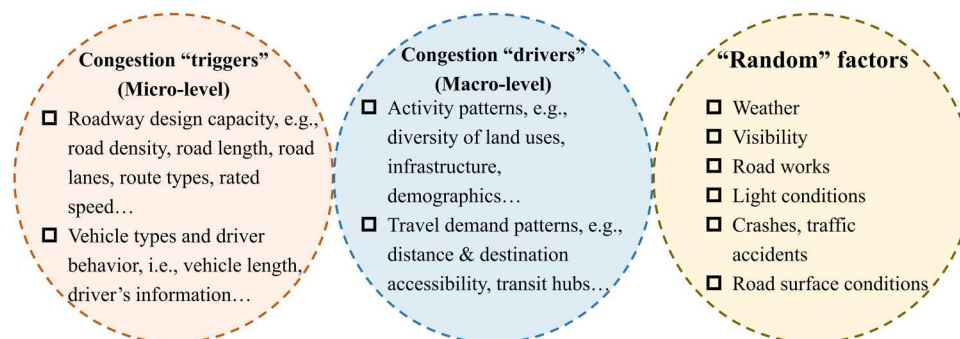
The effects of these factors have been studied by using regression methods (Shen et al. 2020; Lee, Lee, and Putri 2025) or geographical detectors (Deng et al. 2022; Ju et al. 2016). For instance, Zhang et al. (2017a, 2017b) focusing on 'Diversity', models the correlation between traffic congestion and points of interest (POI), and reveals that commercial land negatively affects congestion. Song et al. (2019) adopts the first four 'Ds' to identify which factors influence spatiotemporal congestion patterns. Bao et al. (2023) further combines land use and transportation network data, discovering that congestion in satellite cities associated with commercial land can be offset by public transit. Finally, Pan et al. (2020) and Olayode et al. (2025) integrate road attributes into their analyzes conclude that public commercial POIs, residential POIs, bus routes, bus stops, road lanes, and traffic volumes are the significant contributors to congestion.

Another branch of research emphasized the effects of a sixth 'D' (Demand management) – such as introducing on-demand ridesharing (Li et al. 2022) and congestion pricing (Cook et al. 2025). For example, Qian et al. (2020) and Diao, Kong, and Zhao (2021) demonstrate that the entry of transportation network companies increases congestion in the United States (US). Rahman et al. (2022) uses data from over 100 metropolitan regions in the US and apply a structural equation model (SEM) to show that indirect effects are strong enough to offset direct effects. Huang and Xu (2023) employ a difference-in-differences (DID) approach to evaluate the effect of dockless bike-sharing service (DBS) entry on urban traffic congestion in 98 cities in China, finding that DBS entry reduced congestion by 2.2%. Liang et al. (2023) also employs DID to assess the impact of congestion taxes on traffic congestion, revealing that congestion pricing policy can only slightly mitigate traffic congestion.

The above studies rely solely on correlation, an approach prone to misidentifying relationships (Kamat 2025) due to its failure to account for latent variables. Without testing causal relations, identifying factors as 'key contributors to congestion' based on correlation alone may result in misguided feature selection in policy models, causing policy misallocation.

Few studies employ causal analyses of traffic congestion. The European Conference of Ministers of Transport identifies three categories of causal factors contributing to congestion (Managing urban traffic congestion 2007) (Figure 1): (i) micro-level factors, e.g. conditions on the roadway, or 'congestion triggers'; (ii) macro-level factors, e.g. road usage demand and exogenous factors related to activity patterns and travel demand, or 'road drivers'; (iii) random factors, e.g. weather, visibility, road work, lighting conditions, crashes, special events, etc. Koźlak and Wach (2018) examine nine macro-level factors in Poland using statistical and regression methods. In Contrast, Pi et al. (2019) focuses on micro-level factors to present a visual cause analysis. Yet, to the best of our knowledge, no existing work explores both micro – and micro-factors together.

To fill the gap, we employ multiple causal inference models to identify true causal factors. Based on this, our main contributions lie in: (1) uniquely distinguishing between correlated and causal factors, thereby



**Figure 1.** Three broad categories of causal factors of traffic congestion.

improving model feature selection accuracy (Rotari and Kulahci 2024; Saarela 2024); (2) synthesizing results across multiple causal inference models, thereby enhancing the robustness of our findings (Dormann et al. 2018) and (3) exploring the causal roles of these factors in relation to the perpetually congested roadways (PCRs), a topic that to our knowledge, has not yet been studied.

The remainder of this article is organized as follows. Section 2 describes the data and methodology; Section 3 presents our results; Section 4 discusses the findings; Section 5 concludes this paper and summarizes the key insights.

## 2. Material and methods

### 2.1. Workflow

Figure 2 shows the workflow of this study. First, a multivariable least squares regression (MLSR) is applied to assess the correlation between traffic congestion and built environment factors, comprising nine micro-level factors and 22 macro-level factors. Second, five causal inference models are employed to detect causal relationships. Then, based on the causal factors, the drivers of each spatiotemporal congestion pattern of PCRs are further explored. The final stage involves deriving key insights and policy implications.



**Figure 2.** Workflow of the study.

## 2.2. Data sources

This study uses New York City (NYC) as its case study, due to the availability of a comprehensive congestion dataset (see below). As of 2018, NYC's population was approximately 8.6 million and is projected to increase to 1 million people by 2030 (Solecki and Rosenzweig 2019). This population growth exacerbates its traffic congestion, highlighting the need to devise new solutions.

*Road network data*. These data were extracted from OpenStreetMap, comprising all 100,206 road segments within NYC.

*Built environment data*. Nine micro-level and 22 macro-level factors were obtained from NYC Open Data. According to Ewing and Cervero (2010), the 31 factors can be classified into 'five Ds' (Table 1). Appendix Notes 1 and 2 explain the rationale for excluding the sixth 'D'.

*Average hourly travel speed data and variable descriptions*. The dataset contains 25,068,883 records of average hourly travel speed data for the period 1–31 Dec 2018, provided by Uber Movement. Each record includes recording time, road segment ID, and average speed (Appendix Table 1). Uber Movement also provides free-flow speed data, defined as 15th percentile value of actual speeds of all floating vehicles, sorted in descending order. Traffic congestion is measured by the travel time index (TTI), calculated as the ratio of free-flow speed to average speed (Kong, Yang, and Yang 2015).

## 2.3. Models

Employing multiple models to reduce bias is a well-established practice in the causal inference literature (Imbens and Rubin 2015). The rationale for selecting the five models used in this study is provided in Appendix Note3.

**Table 1.** Definition and calculation methods of the 31 representative factors.

| Five 'Ds' and definitions | Representative variables | Calculation methods |
|---|---|---|
| Diversity (The degree of land use mixture within a defined area) | Proportion of listed POIs, including commercial, residential, health services, social services, cultural services, education, recreational, government, transportation, public safety, and waterbody. | Number of POIs/Area of statistical unit (i.e. the buffer with 700m buffer around each road segment; Song et al. 2019) |
| Density (Concentration of population or physical structures within a given area) | Population density | $\sum_{i \in J} P_i / Area_J$, where $P$ is the population in the 100m pixel $i$ located within statistical unit $J$ $Area_J$ is the area of $J$. |
| | Building density | $\sum_{i \in J} BA_i / Area_J$, where $BA_i$ is footprint area of the building. |
| | Ramp density | $\sum_{i \in J} LR_i / Area_J$, where $LR_i$ is the length of ramp. |
| | Road density | $\sum_{i \in J} LRD_i / Area_J$, $LRD_n$ is the length of the road. |
| | Parking lot density | $\sum_{i \in J} PA_i / Area_J$, where $PA_i$ is the area of parking lot $i$. |
| | Pedestrian zone density | $\sum_{i \in J} PZA_i / Area_J$, where $PZA_i$ is the area of pedestrian zone $i$. |
| | Bus stop shelter density | $\sum_{i \in J} BS_i / Area_J$, where $BS_i$ is the number of bus stop shelters with $J$. |
| | Subway station density | $\sum_{i \in J} SS_i / Area_J$, $SS_n$ is the number of subway stations within $J$. |
| Design (Configuration and physical attributes of transportation infrastructure) | Number of road lanes | Number |
| | Length of road segments | Number |
| | Ratio of road length to width | Ratio = road length / road width |
| | Number of traffic lights | Count within statistical unit |
| | Number of street lights | Count within statistical unit |
| | Length of bike routes | Cumulative length within statistical unit |
| | Length of truck routes | Cumulative length within statistical unit |
| | Length of sidewalk routes | Cumulative length within statistical unit |
| Distance to Transit & Destination Accessibility (Proximity to transit hubs and access to public transportation nodes) | Distance to the nearest bridge | Euclidean distance (in km) from road segment centroid to each respective facility. |
| | Distance to the nearest bus stop | |
| | Distance to the nearest subway | |
| | Distance to the nearest airport | |

### 2.3.1. Correlation model

The MLSR was applied to analyze correlations due to its simplicity and wide use in the literature (Bao et al. 2022; Bao et al. 2023; Koźlak and Wach 2018; Shen et al. 2020). The model is specified as

$$TTI_{weighted} = \beta_0 + \beta_i X_1 + \beta_2 X_2 + \ldots + \beta_J X_J, \tag{1}$$

where $TTI_{weighted} = [\overline{TTI}_1, \ldots, \overline{TTI}_N]^T$ represents the overall congestion levels of N road segments, and $\overline{TTI}_i$ is the entropy-weighted average of TTIs over all days. An example of calculating $\overline{TTI}_i$ is provided in Appendix Note 4. $X_j \in R^{N \times 1}$ denotes the $j_{th}$ factor ($J = 31$), and $\beta_j$ are the parameters to be estimated.

To address multicollinearity and high inter-variable correlation, we additionally employed variance inflation factor (VIF) and Pearson correlation coefficient (Shrestha 2020; Tay 2017). A VIF $\in [0, 5]$ typically indicates the absence of multicollinearity.

### 2.3.2. Causality models

Granger causality (Shojaie and Fox 2022) is a prediction-based method. If the historical data of a particular built environment factor improves the prediction accuracy of congestion level, then that factor is said to 'G-causes' congestion. The method involves three stages: first, a baseline model is constructed to predict future congestion using only historical congestion data. Second, an enhanced model is constructed that incorporates historical data on the built environment. Finally, the predictive performance of the baseline model and enhanced model is compared. If the enhanced model performs better, it indicates that the factor 'G-causes' congestion. Traditional Granger causality assumes linear relationships. To address this limitation, we adopt an improved model based on transfer entropy (TE), and calculate the net information outflow (Wiener 1956) to quantify causality:

$$\widehat{TE}_{factor \to congestion} = TE_{factor \to congestion} - TE_{congestion \to factor}, \tag{2}$$

if $\widehat{TE}_{factor \to congestion} \geq 0$, this implies that the factor 'G-causes' congestion.

SEM (Golob 2003) distinguishes between direct and indirect effects of the 'five Ds' and their composite variables (i.e. XX1-XX5). In contrast to Granger causality, SEM captures the pathways through which the 'five Ds' influence traffic congestion by incorporating latent variables. The model is expressed as follows:

$$TTI \sim XX1 + XX2 + XX3 + XX4 + XX5, \tag{3}$$

$$XX1 =\sim XX5, \ XX2 =\sim XX5, XX3 =\sim XX5, XX4 =\sim XX5, \tag{4}$$

$$XX1 \sim\sim XX2, XX2 \sim\sim XX3, XX3 \sim\sim XX4, XX4 \sim\sim XX1. \tag{5}$$

Here, the symbols ($\sim$, $=\sim$, $\sim\sim$) denote regression, latent variable definition and variance, respectively. XX1-XX4 are calculated as the weighted averages of components of the 'five Ds', while XX5 represents the average XX1 – XX4.

Causal forest outperforms the previous two models because it can identify the congestion factors at a granular level, by estimating the heterogeneous treatment effects (HTEs) (Wager and Athey 2018). For example, while traditional modes may indicate how to alleviate congestion at the citywide level, causal forest provides evidence for specific areas. It functions like an intelligent diagnostic system, capable of identifying different 'treatment plans' for different road segments-an approach useful for policymakers focused on localized interventions.

Causal impact (Brodersen et al. 2015) differs from causal forest by comparing the congestion levels under conditions where a specific built environment factors is considered versus where is it. This method creates a treatment group and a control group, and follows a three-stage procedure: first, a Bayesian structural time series model is constructed to learn the effect pattern of the given factor on congestion; second, the model is used to generate a 'synthetic control' time series; and third, the observed effects are compared between the treatment and control groups, capturing both pointwise and cumulative impacts.

Convergent cross mapping (CCM) (Tsonis et al. 2018) measures the causal relationship between traffic congestion Y and a built environment factor X based on their short-term temporal dynamics. By constructing their respective manifolds, $M_Y$ and $M_X$, under the embedded dimension $E$, the $E + 1$ nearest points to

$X(t)$ in $M_X$ and to $Y(t)$ in $M_Y$ can be used to estimate the system as follows:

$$\hat{Y}(t)|M_X = \sum w_i Y(t_i), \ i = 1, \ldots, \ E+1, \tag{6}$$

$$\hat{X}(t)|M_Y = \sum w_i X(t_i), \ i = 1, \ldots, \ E+1, \tag{7}$$

where $w_i$ is the weight assigned to each neighboring point. The CCM correlation coefficient is then calculated as

$$r_{CCM} = \frac{\sum_{i=1}^{L}(X(i) - \overline{X(i)})(\hat{X}(i)|M_Y - \overline{\hat{X}(i)|M_Y})}{\sqrt{\sum_{i=1}^{L}(X(i) - \overline{X(i)})^2 \sum_{i=1}^{L}(\hat{X}(i)|M_Y - \overline{\hat{X}(i)|M_Y})^2}}. \tag{8}$$

As the time series L increases, if $r_{CCM}$ converges to a stable value, it indicates a causal relationship from $X$ to $Y$.

The five causal models offer complementary approaches for distinguishing relationships between the built environment and traffic congestion. Granger causality acts as a 'time-series prediction', assessing whether historical built environment data can predict future congestion – a straightforward approach but sensitive to the choice of time window. SEM serves as a 'roadmap architect', quantifying both direct effects (e.g. the impact of road density on traffic congestion) and indirect pathways (e.g. the influence of bus stops on congestion through transit use). Causal Forest functions like a 'precision tool', identifying location-specific effects, through is computationally intensive. Causal Impact is well-suited for evaluating short-term effects using synthetic controls, yet is vulnerable to control group selection bias. CCM operates as a 'long-term observer', capturing delayed and nonlinear responses (e.g. the influence of water bodies may take years to manifest in traffic patterns). For comprehensive analysis, we integrate these models to uncover how the built environment influences congestion through multiple causal pathways.

## 2.4. Perpetually congested roadway (PCRs) analysis

Congestion in PCRs is more intricate than general congestion. In this study, road segments with TTI of at least 1.5 for all hours of the day are defined as PCRs. This threshold is selected based on the classification proposed by Kong, Yang, and Yang (2015), where TTI≤1.5 indicates smooth or free-flowing traffic (see Appendix Table 2).

### 2.4.1. Congestion patterns

Agglomerative hierarchical clustering (Oti and Olusola 2024) is employed to identify spatiotemporal patterns within PCRs based on the TTI. Four linkage methods are commonly used to measure distance: single, complete, average, and Ward linkage. Ward linkage is selected for this study, and the distance between two clusters $(C_i, C_j)$ is defined as

$$D_{ward}(C_i, C_j) = \frac{1}{|C_i| \cdot |C_j|} \sum_{x_i \in C_i, x_j \in C_j} (dist(x_i, x_j))^2, \tag{9}$$

where $dist(x_i, x_j)$ denotes the Euclidean distance between samples $x_i$ and $x_j$. $|C_i|$ and $|C_j|$ represent the number of samples in clusters $C_i$ and $C_j$. Subsequently, all PCRs are clustered using k-means, with the weighted feature matrix, expressed as:

$$\hat{F} = F*diag(w_1, w_2, \ldots, w_{11}), \tag{10}$$

where $F \in R^{N*11}$ is the feature matrix. $N$ is the number of road segments, $w_j$ ($j = 1, 2, \ldots 11$) denotes the weight for the $j^{th}$ feature, and $\hat{F} \in R^{N*11}$ represents the resulting weighted feature matrix.

### 2.4.2. Mantel test

The Mantel test (Somers and Jackson 2022) is conducted to evaluate the association between the TTI matrices of the congestion patterns and the feature matrix $F$. The Mantel test is a statistical method used

for evaluating the correlation between two data matrices and to test for autocorrelation. We chose the Mantel test because, unlike conventional correlation coefficients that measure relationships between pairs, it is specifically designed to evaluate correlations between entire matrices. Unlike the Pearson and Spearman coefficients, the Mantel test uses a permutation-based approach to randomly rearranging the data. Each permutation is calculated only once. The Mantel p correlation ranges from $-1$ to 1, where $-1$ indicates a negative correlation, 0 indicates no correlation, and 1 indicates a perfect positive correlation.

# 3. Results

## 3.1. Single-model results and validity

### 3.1.1. MLSR

(1) Model results. As shown in Table 2, six factors do not exert a statistically significant influence ($P > 0.05$) and three factors exhibit multicollinearity (VIF > 6). Combined with the Pearson correlation coefficients (Appendix Figures 1 and 2), only 19 factors are found to be significantly correlated with the TTI (Appendix Table 3). The $R^2$ value from the MLSR is relatively low ($R^2 = 0.150$). Durbin-Watson statistics and Hypothesis tests confirm that our data satisfy the assumptions of the regression model except for linearity. To further validate these results, we applied a non-linear model (random forest), which achieved a substantially better performance ($R^2 = 0.565$; Table 4 and Figure 3 in Appendix).

(2) Model validity. As shown in Table 2, the values of Mean Square Error (MSE), Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and Mean Absolute Percentage Error (MAPE) are all close to zero, indicating that the model exhibits strong predictive performance.
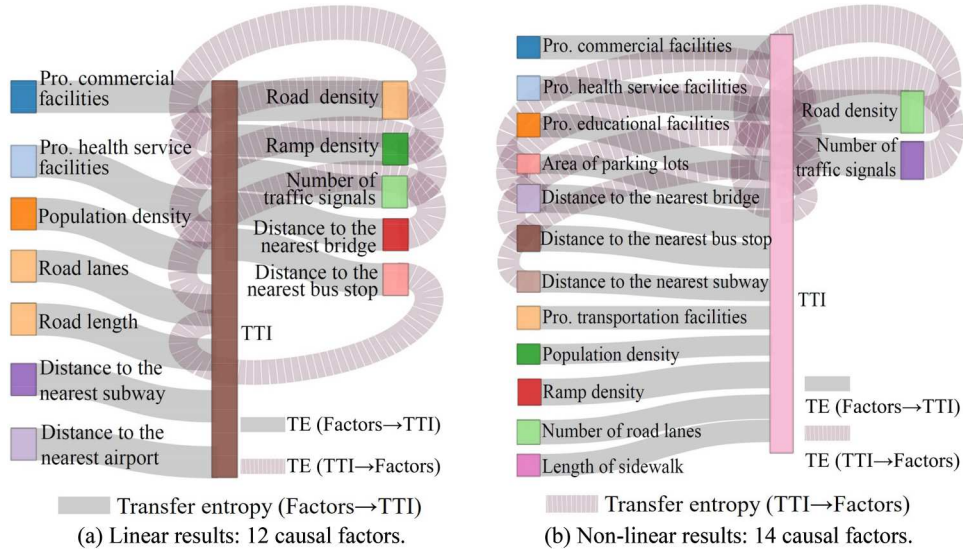
### 3.1.2. Granger causality

(1) Model results. Figure 3 visualizes the results of both the linear and nonlinear models, with corresponding numerical values presented in Table 3. The analysis reveals that the proportion of commercial and

**Table 2.** MLSR analysis between 31 built environment factors and traffic congestion.

| Five 'Ds' | Representative variables | Coefficient | t-Value | P-value | VIF |
|---|---|---|---|---|---|
| Diversity | Proportion of commercial facilities | 0.078 | 20.194 | *** | 4.155 |
| | Proportion of residential facilities | 0.131 | 2.762 | *** | 17.149 |
| | Proportion of health service facilities | 0.019 | 7.958 | *** | 2.034 |
| | Proportion of social service facilities | 0.026 | 1.061 | 0.289 | 3.135 |
| | Proportion of cultural service facilities | −0.005 | −1.638 | 0.101 | 18.352 |
| | Proportion of educational facilities | 0.011 | 7.103 | *** | 2.210 |
| | Proportion of recreational facilities | 0.024 | 6.362 | *** | 2.269 |
| | Proportion of government facilities | 0.020 | 6.762 | *** | 1.553 |
| | Proportion of transportation facilities | −0.025 | −5.928 | *** | 1.743 |
| | Proportion of public safety facilities | −0.002 | −0.808 | 0.419 | 1.305 |
| | Proportion of waterbody | −0.001 | −0.241 | 0.810 | 1.759 |
| Density | Population density | 0.029 | 5.313 | *** | 1.223 |
| | Building density | 0.012 | 2.150 | ** | 1.811 |
| | Road density | 0.008 | 3.945 | *** | 1.760 |
| | Ramp density | 0.044 | 16.176 | *** | 1.127 |
| | Density of parking lots | 0.028 | 3.880 | *** | 1.043 |
| | Density of pedestrian zones | 0.006 | 0.818 | 0.413 | 1.115 |
| | Density of bus stop shelters | −0.025 | −10.896 | *** | 2.576 |
| | Density of subway stations | 0.041 | 9.766 | *** | 3.357 |
| Design | Number of road lanes | 0.023 | 6.767 | *** | 5.228 |
| | Length of road segments | −0.036 | −2.266 | ** | 1.656 |
| | Ratio of road segment length to width | 0.184 | 5.025 | *** | 6.609 |
| | Number of traffic lights | 0.033 | 13.136 | *** | 2.085 |
| | Number of street lights | 0.024 | 5.948 | *** | 2.712 |
| | Length of bike routes | 0.024 | 8.095 | *** | 1.065 |
| | Length of truck routes | −0.008 | −1.588 | 0.112 | 1.053 |
| | Length of sidewalk routes | 0.014 | 3.323 | *** | 1.037 |
| Distance & Destination Accessibility | Distance to the nearest bridge | −0.012 | −5.029 | *** | 1.316 |
| | Distance to the nearest bus stop | −0.060 | −16.759 | *** | 1.339 |
| | Distance to the nearest subway | −0.037 | −19.711 | *** | 1.616 |
| | Distance to the nearest airport | −0.022 | −13.022 | *** | 1.236 |

MSE = 0.003, RMSE = 0.058, MAE = 0.042, MAPE $= 1 \times 10^{-6}$, $R^2 = 0.150$

**Figure 3.** Granger causality results. Factors connected by loops represent bidirectional causality. Only road density and the number of traffic lights show robust bidirectional causality with the TTI, as reflected by their larger node sizes.

health service facilities, population density, number of road lanes and traffic lights, road density, distance to the nearest bridge and bus stop are statistically significant causal factors ($P_{X \to Y} < 0.05$). Among these, road density ($TE_{Y \to X} = 0.177$) and the number of traffic lights ($TE_{Y \to X} = 0.116$) exhibit the strongest bidirectional causality with the TTI, suggesting these elements not only influence but are also influenced by traffic conditions in a feedback loop.

(2) Statistical test. As shown in Table 3, a $P$-value less than 0.05 indicates a statistically significant causal effect.

(3) Parameter setting. Figure 4 illustrates how the $p$-value varies with lag length, based on the Chi-squared and the Likelihood ratio tests. $p$-value less than 0.05 indicates reasonable lag. The optimal settings identified across the tests are 100, 20, 95, 115, 125, 130, 95, 40, 75, 130, 75, 60, 80, 85, 140, and 140.

### 3.1.3. SEM

(1) Model results. In Figure 5a, 14 factors emerge as statistically significant causal variables, indicated by solid arrows. XX4 stands out as particularly influential because all of its component variables (X16-X19) achieve statistical significance, suggesting its causal effect. Figure 5b presents the overall effects of XX1-XX4, with specific effect paths detailed in Appendix Table 5. Intuitively, XX4 (Distance and Destination accessibility) shows the strongest direct effect, suggesting it may represent a fundamental determinant of traffic congestion. XX1 (Diversity) demonstrates the strongest indirect effect, highlighting how urban characteristics influence traffic through secondary pathways. These findings underscore the importance of both direct and indirect of urban factors' effects when designing policy inventions.
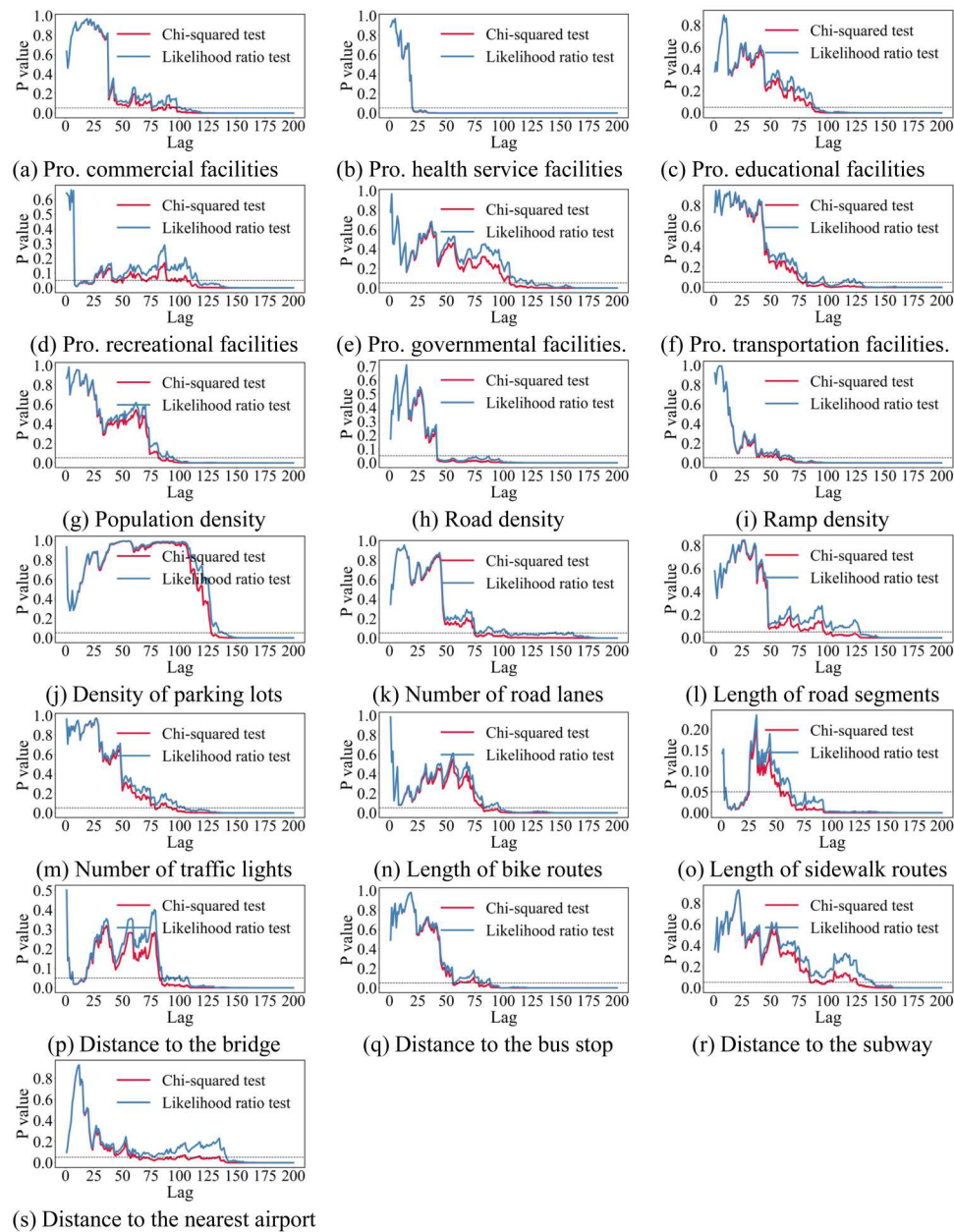
(2) Validity test. The model's fit is evaluated using standard goodness-of-fit statistics (Rahman et al. 2022). As shown in Table 4, the model yields a statistically significant CMIN score of 63.149. The GFI exceeds the recommended threshold of 0.95, indicating a strong model fit. The CFI, NFI, and IFI with values of 0.905, 0.094, and 0.907, respectively-also support the model's adequacy.

### 3.1.4. Causal forest

(1) Model results. In Figure 6, SHapley Additive exPlanations (SHAP; Albini et al. 2022) are used to visualize feature changes. SHAP values quantify feature effects for individual road segment, while the final effect is aggregated across all segments. Taking X16 as an example, the SHAP values of most road segments are greater than zero, indicating that its overall positive influence outweighs any negative effects. Similar analysis applies to other factors. In Figure 7a, X8 and X16 show strong feature importance. In Figure 7b, X16, X7, X15, X17, and X12 show strong importance, as indicated by the feature splitting order. These correspond to

**Table 3.** Statistical test of granger causality.

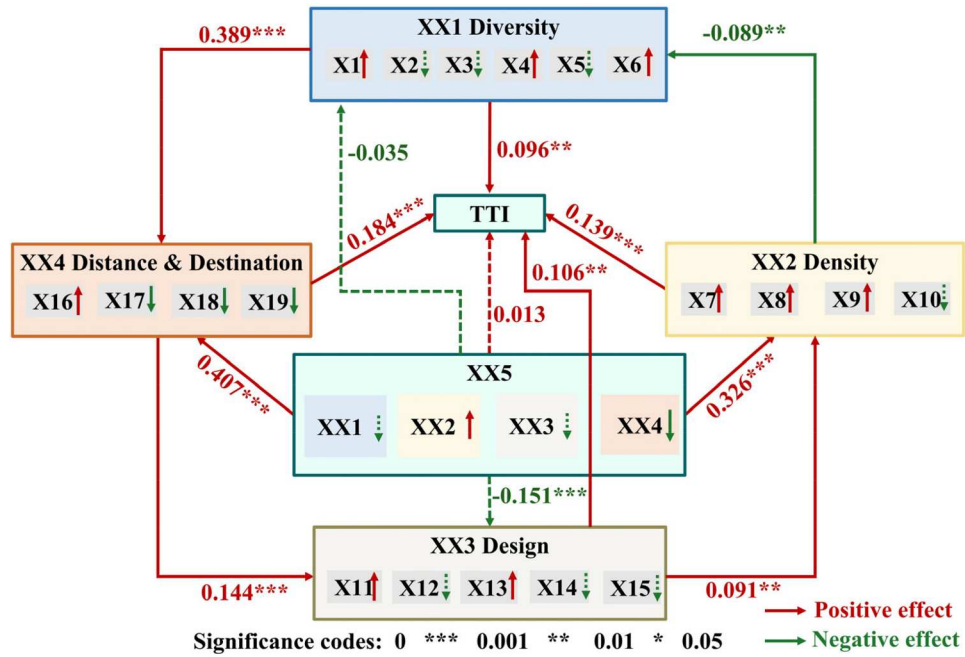| Variable | Linear transfer entropy based-method | | | Nonlinear transfer entropy based-method | | |
|---|---|---|---|---|---|---|
| | P-value ($P_{X\rightarrow Y}$, $P_{Y\rightarrow X}$) | Z-score ($Z_{X\rightarrow Y}$, $Z_{Y\rightarrow X}$) | TE ($TE_{X\rightarrow Y}$, $TE_{Y\rightarrow X}$) | P-value ($P_{X\rightarrow Y}$, $P_{Y\rightarrow X}$) | Z-score ($Z_{X\rightarrow Y}$, $Z_{Y\rightarrow X}$) | TE ($TE_{X\rightarrow Y}$, $TE_{Y\rightarrow X}$) |
| X1 | (0.000, 0.067) | (4.401, 2.249) | (0.001, 0.001) | (0.000, 0.733) | (3.590, −0.785) | (0.073, 0.065) |
| X2 | (0.000, 0.053) | (1.756, −0.440) | (0.001, 0.001) | (0.000, 0.200) | (2.282, 1.165) | (0.081, 0.068) |
| X3 | (0.067, 0.133) | (1.686, 0.816) | (0.000, 0.001) | (0.000, 0.667) | (2.367, −0.576) | (0.080, 0.072) |
| X4 | (0.333, 0.267) | (0.185, 0.142) | (0.001, 0.001) | (0.133, 0.200) | (1.220, 0.493) | (0.093, 0.096) |
| X5 | (0.133, 0.400) | (0.560, 0.130) | (0.000, 0.001) | (0.400, 0.400) | (0.341, 0.358) | (0.094, 0.103) |
| X6 | (0.067, 0.200) | (2.239, 0.689) | (0.001, 0.001) | (0.000, 0.333) | (1.851, 1.001) | (0.077, 0.066) |
| X7 | (0.000, 0.533) | (3.076, −0.392) | (0.001, 0.001) | (0.000, 0.333) | (1.798, −0.034) | (0.067, 0.059) |
| X8 | (0.000, 0.033) | (3.790, −0.770) | (0.001, 0.001) | (0.000, 0.033) | (1.66, −1.393) | (0.096, 0.117) |
| X9 | (0.000, 0.047) | (3.123, 1.367) | (0.001, 0.001) | (0.000, 0.067) | (2.028, −0.471) | (0.088, 0.100) |
| X10 | (0.133, 0.867) | (0.913, −0.654) | (0.000, 0.001) | (0.000, 0.000) | (2.342, 2.160) | (0.068, 0.055) |
| X11 | (0.000, 0.667) | (3.432, −0.694) | (0.001, 0.001) | (0.400, 0.200) | (0.071, 0.472) | (0.082, 0.088) |
| X12 | (0.000, 0.533) | (7.790, −0.436) | (0.001, 0.001) | (0.000, 0.533) | (1.245, −0.225) | (0.089, 0.091) |
| X13 | (0.000, 0.033) | (5.436, −0.517) | (0.001, 0.001) | (0.000, 0.000) | (2.499, −3.541) | (0.103, 0.116) |
| X14 | (0.500, 0.600) | (2.600, 21.910) | (0.001, 0.000) | (0.067, 0.467) | (1.195, 0.285) | (0.090, 0.099) |
| X15 | (0.267, 0.100) | (1.083, 3.573) | (0.000, 0.000) | (0.000, 0.667) | (2.056, −0.535) | (0.074, 0.066) |
| X16 | (0.000, 0.000) | (3.881, 2.216) | (0.001, 0.001) | (0.000, 0.000) | (1.791, 1.890) | (0.090, 0.093) |
| X17 | (0.000, 0.049) | (2.439, −0.858) | (0.001, 0.001) | (0.000, 0.000) | (1.296, −1.181) | (0.086, 0.079) |
| X18 | (0.000, 0.467) | (3.352, −0.176) | (0.001, 0.000) | (0.000, 0.000) | (1.864, −0.548) | (0.059, 0.067) |
| X19 | (0.000, 0.133) | (2.922, 1.528) | (0.001, 0.001) | (0.267, 0.933) | (0.516, −1.260) | (0.091, 0.096) |

**Figure 4.** Lag selection of the model. This refers to the number of prior observation points included in the model.
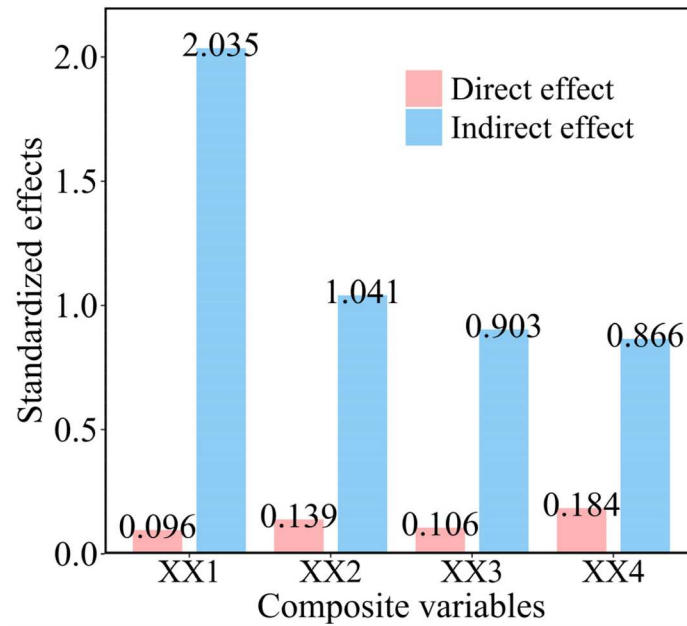
**Table 4.** Evaluation of fitting goodness.

| Index | Recommended score | Results of our model |
|---|---|---|
| Chi-square (CMIN) | A lower value indicates a better fit. | 63.027 |
| Goodness-of-fit index (GFI) | [0,1], 1 indicates perfect fit. | 0.973 |
| Comparative fit index (CFI) | [0,1], 1 indicates perfect fit. | 0.905 |
| Normal fit index (NFI) | [0,1], 1 indicates perfect fit. | 0.904 |
| Incremental fit index (IFI) | [0,1], 1 indicates perfect fit. | 0.907 |

distance to the nearest bridge, population density, sidewalk route length, distance to the nearest bus stop, and road length-highlighting their strong causal relationship with the TTI.

(2) Model robustness. Figure 8 presents the comparison between the estimated treatment effects generated by the model and the true treatment effects for the six factors. All estimated effect curves exhibit trend consistent with the true curves, indicating strong predictive performance and suggesting good model generalization.

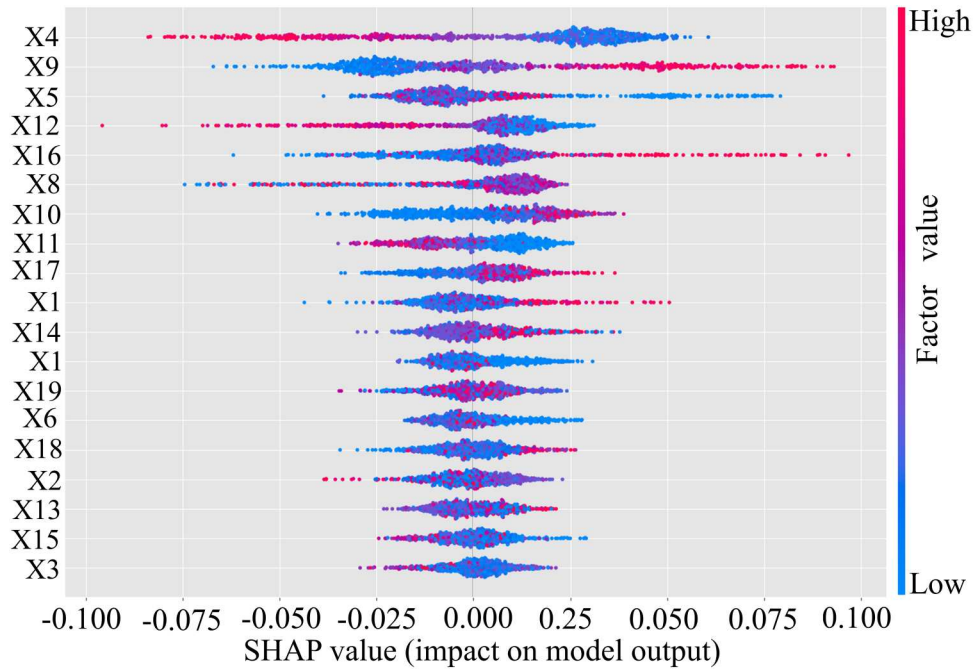(a) Effects and significance of X1-X19, and XX1-XX5.



(b) Direct and Indirect effects of XX1-XX5.

**Figure 5.** SEM results. Four factors in XX4 are all significant, which show the most robust causality.

### 3.1.5. Causal impact

(1) Model results. As shown in Figure 9, ten factors show causal relationships with the TTI in this model. The dataset is partitioned with a training to test ratio of 2:1. When the sample size exceeds 200, both the pointwise and the cumulative plots show an upward trend, indicating a positive causal effect on congestion. Conversely, a downward trend suggests a negative causal effect.

(2) Statistical test. In Table 5, the posterior causal probability reflects the likelihood that a variable truly influences the TTI. Compared to the $p$-value, it offers a more intuitive measure of causal certainty. We rigorously employ both $p$-value and poster probability to ensure a robust result. The results identify ten factors with $p$-value less than 0.05 and posterior probability exceeding 95%, thereby providing strong evidence of a causal effect.

**Figure 6.** Causal forest results. The effects of X1-X19 on individual read segment. The larger range of the points distribution, the stronger impacts.

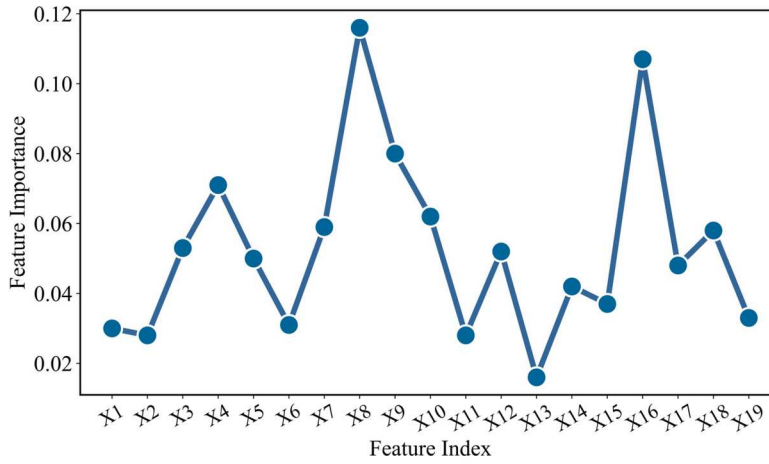**Table 5.** Statistical test of causal impact.

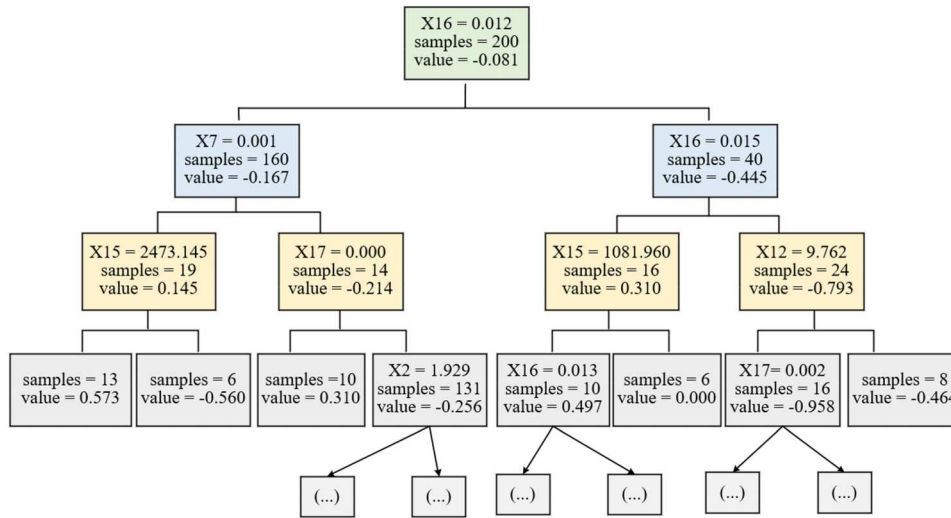| Variable | Actual value -average | Prediction value- average | Actual value- cumulative | Prediction value- cumulative | Absolute effect- average | Absolute effect- cumulative | Relative effect | P- value | Posterior probability of a causal effect |
|---|---|---|---|---|---|---|---|---|---|
| X1 | 2.80 | 2.00 | 276.80 | 197.00 | 0.80 | 80.00 | 49% | 0.011 | 98.90% |
| X2 | 2.70 | 2.00 | 266.60 | 196.00 | 0.71 | 70.52 | 44% | 0.056 | 94.00% |
| X3 | 2.90 | 1.90 | 292.50 | 192.90 | 1.00 | 100.00 | 69% | 0.039 | 96.07% |
| X4 | 2.90 | 3.20 | 285.10 | 322.70 | −0.83 | −37.63 | −11% | 0.083 | 92.00% |
| X5 | 2.60 | 2.00 | 258.60 | 195.0 | 0.63 | 63.34 | 40% | 0.069 | 93.00% |
| X6 | 2.80 | 3.50 | 278.30 | 347.70 | −0.69 | −69.39 | −18% | 0.083 | 92.00% |
| X7 | 3.10 | 2.50 | 312.60 | 253.80 | 0.59 | 58.85 | 26% | 0.049 | 95.10% |
| X8 | 2.90 | 2.90 | 293.90 | 291.20 | 0.03 | 2.629 | 1.60% | 0.439 | 56.00% |
| X9 | 2.80 | 3.80 | 282.70 | 382.30 | −1.00 | −100.00 | −25% | 0.006 | 99.40% |
| X10 | 2.70 | 2.40 | 269.30 | 239.50 | 0.30 | 29.80 | 30% | 0.336 | 66.00% |
| X11 | 2.70 | 3.20 | 274.00 | 317.30 | −0.43 | −43.31 | −13% | 0.036 | 96.40% |
| X12 | 2.80 | 3.20 | 284.90 | 316.50 | −0.32 | −31.65 | −8.4% | 0.388 | 61.00% |
| X13 | 2.90 | 4.10 | 294.50 | 408.20 | −1.10 | −113.70 | −27% | 0.007 | 99.30% |
| X14 | 2.90 | 2.90 | 293.80 | 286.70 | 0.071 | 7.09 | 3.40% | 0.386 | 61.00% |
| X15 | 3.00 | 2.80 | 299.00 | 284.20 | 0.15 | 14.57 | 8.10% | 0.370 | 63.00% |
| X16 | 3.30 | 2.50 | 326.70 | 246.70 | 0.80 | 80.00 | 35% | 0.009 | 99.05% |
| X17 | 2.80 | 2.10 | 277.00 | 209.90 | 0.67 | 67.07 | 36% | 0.049 | 95.10% |
| X18 | 2.90 | 3.70 | 292.30 | 370.50 | −0.78 | −78.21 | −20% | 0.045 | 95.50% |
| X19 | 2.60 | 3.30 | 267.80 | 328.40 | −0.71 | −70.59 | −21% | 0.011 | 98.90% |

### 3.1.6. CCM

(1) Model results. In Figure 10, an increase in the correlation coefficient indicating a causal relationship. Nine factors exhibit a unidirectional causal effect (i.e. x causes y) with the TTI. Notably, the number of traffic lights exhibit bidirectional causality with the TTI, as evidence by both the x causing y and y causing x effects displaying a rising trend (Figure 10e). The analysis reveals a feedback loop where increased traffic congestion drives additional traffic signal installation, which subsequently alters traffic flow until reaching an observable stable state in the curves. This mutually reinforcing pattern in both directions suggest a complex feedback loop that has significant implications for transportation planning – traffic signal optimization should be viewed as ongoing, dynamic process rather than a one-time solution.

(2) Parameter settings. Figure 11 illustrates how prediction performance varies with the embedding dimensions. The X-axis value corresponding to the peak of each curve represents the optimal embedding

(a) Feature importance of the 19 factors. X8 and X16 show strong importance.



(b) Causal tree. The splitting order is X16, X7, X15, X17, and X12.

**Figure 7.** Feature selection for causal forest model. The variable with high feature importance and a high splitting order is dominant causal factor.
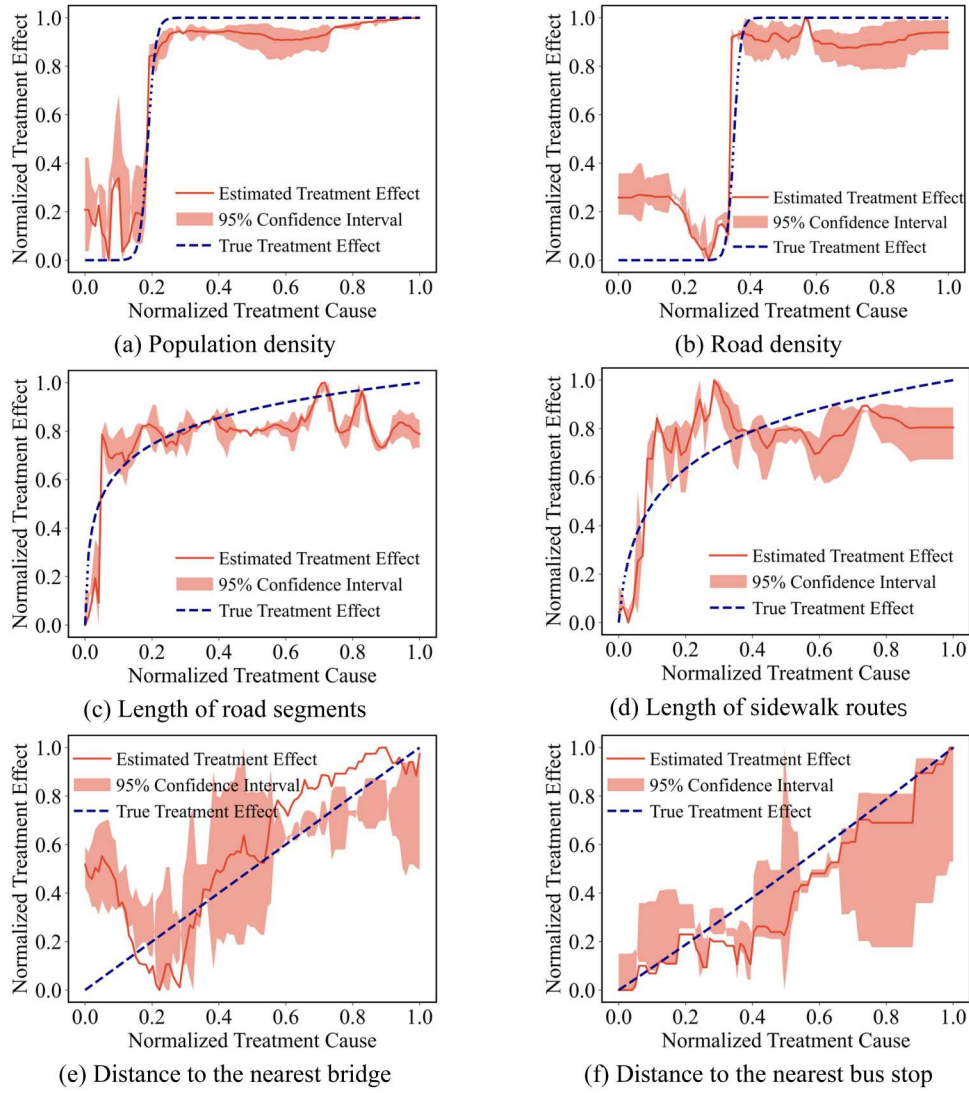
dimension. For example, in the first panel, when exploring the causality from x to the TTI, the optimal embedding dimension of x is 10, while the optimal dimension for TTI is 6.

## 3.2. Overall results

### 3.2.1. Integrated findings

The synthesized results, as illustrated in Figure 12, reveal a hierarchy of causal factors based on their validation across five distinct models. The robustness of each factor was assessed by the number of models in which it demonstrated statistical significance, with a predefined threshold requiring validation in at least three models (i.e. demonstrating minimum 60% causal support). First, the distance to the nearest bridge and bus stop emerges as the most robust causal factors, validated by all five models. This strongly suggests that transportation accessibility plays a fundamental role in the observed phenomenon. Second, several factors showed high but slightly less consistent support, validated in four models, including the proportion of commercial facilities, population density, number of road lanes and traffic lights, and distance to the nearest airport. This indicates the importance of urban infrastructure and demographic characteristics. Then,
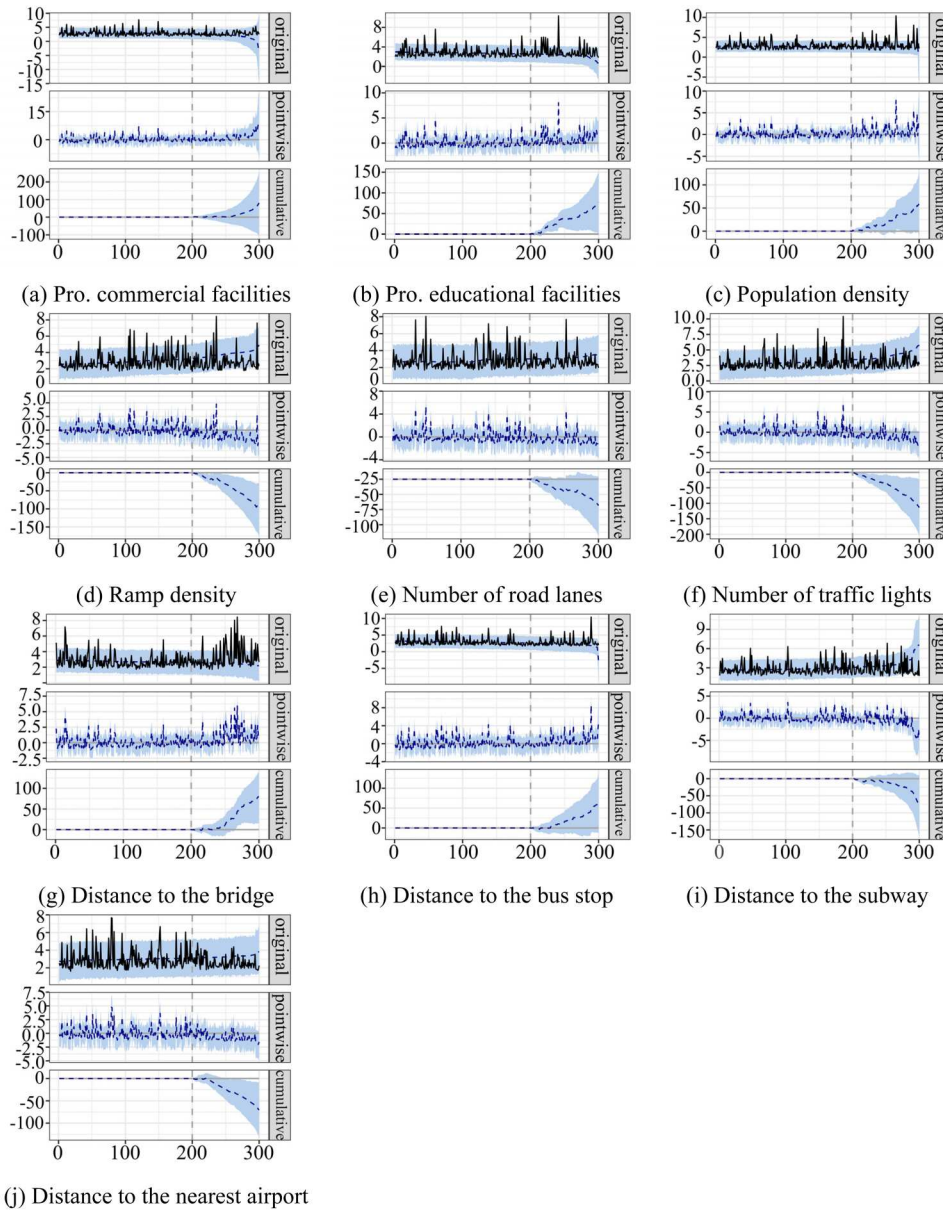
**Figure 8.** Illustration of estimated treatment effect and true treatment effect of causal forest model.

educational facilities, road density, ramp density, and distance to the nearest subway pass validation in three model tests. The remaining factors, validated by only one or two models, demonstrated insufficient consistency for robust conclusions. Overall, this analysis identified 11 robust causal factors, with the findings particularly highlighting transportation accessibility as key determinants. It also suggests that less consistent factors may require further investigation to understand their potential context-dependent effects or possible interactions with other variables.

### 3.2.2. Sensitivity to weekday-weekend-holidays

To assess the robustness of our causal findings, we conducted a comparative analysis using the average treatment effect (ATE) (Li 2020) across four temporal groups: one-month (baseline), weekday, weekend, and holiday (Figure 13). First, the analysis revealed consistent causal rankings for infrastructure-related factors, such as distance to bridges and transit hubs, number of traffic lights and of road lanes, which consistently ranked among the top-six in all subgroups. This confirms their temporal stability. However, the weekend and holiday groups exhibited notable shifts in specific land use categories. Compared to the baseline group, the ATE for commercial facilities and recreational facilities increased by a 48.6% and 60.0%, respectively, while educational facilities showed a 64.8% decrease. This is consistent with travel behavior theory (Zhu et al. 2019). Third, weekday results displayed a trend aligned with the full-month pattern, suggesting that weekday data effectively captures representative congestion dynamics.

(a) Pro. commercial facilities  (b) Pro. educational facilities  (c) Population density

(d) Ramp density  (e) Number of road lanes  (f) Number of traffic lights

(g) Distance to the bridge  (h) Distance to the bus stop  (i) Distance to the subway
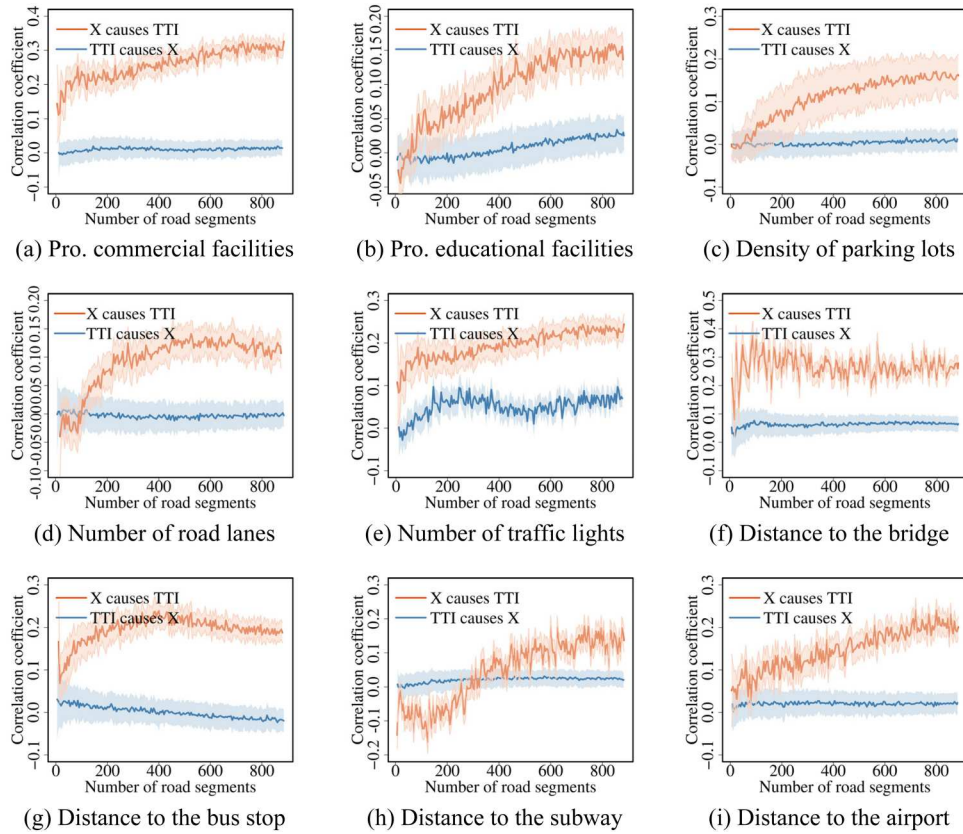
(j) Distance to the nearest airport

**Figure 9.** Causal impact results. For each figure, the trend of pointwise and cumulative effects represents the variation of causal effects on traffic congestion. Five factors show positive upward trend and the other five factors show negative downward trend.

In summary, our findings indicate while core transit infrastructures effects remain robust across temporal contexts, land use impacts are temporally sensitive. Urban planners can confidently rely on infrastructure-related findings but should adjust for behavioral shifts in land use effects during non-workday periods.

### 3.3. Perpetually congested roadway

#### 3.3.1. Identified spatiotemporal congestion patterns

Our procedure identifies 889 segments classified as PCRs. The hierarchical clustering procedure identifies four distinct spatiotemporal congestion patterns (Appendix Figures 4), which are subsequently grouped into two higher-level clusters using k-means based on the 11 weighted causal factors. The silhouette coefficient (Rousseeuw 1987) and elbow method (Liu and Deng 2020) both indicate that two clusters are the optimal solution, as demonstrated in Appendix Figure 5. These clusters are interpreted as main urban arterials

**Figure 10.** CCM results. The trends represent the variation of causal effects. Number of traffic lights shows bidirectional causality, as its two lines both show upward trend.

(cluster 1) and city expressways (cluster 2), respectively, with their spatial distributions illustrated in Appendix Figure 6. The corresponding statistical summaries of these clusters are comprehensively presented in Figure 14. This two-step clustering approach successfully reveals the two main types of congestion patterns in city road networks.
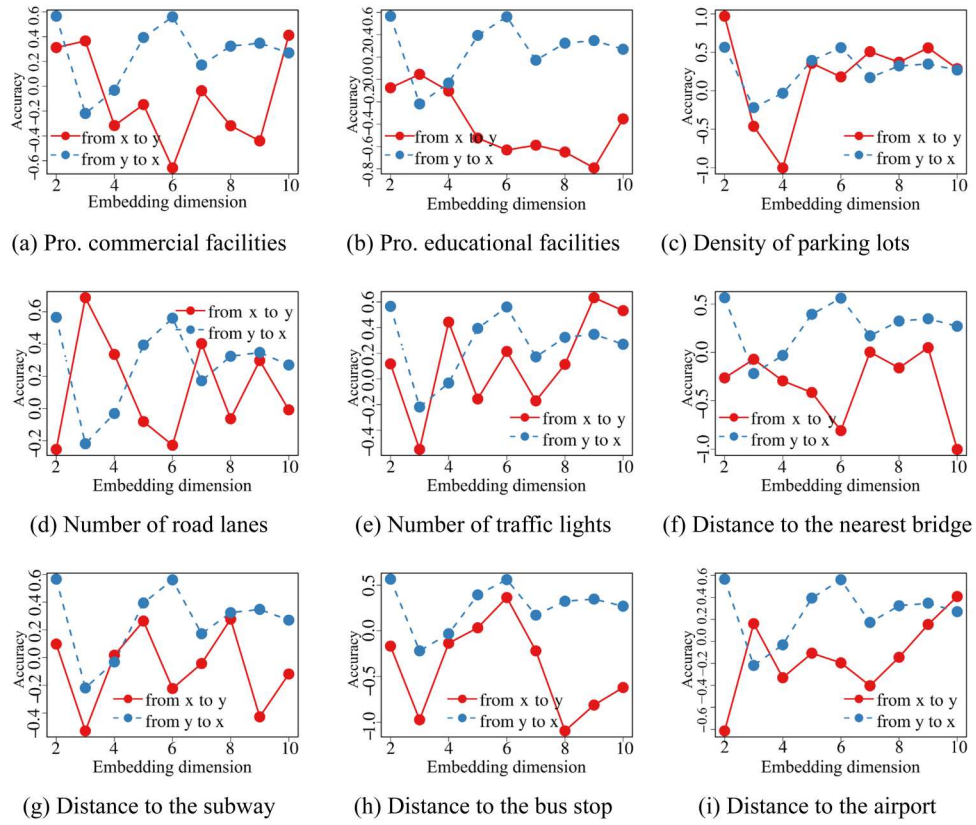
### 3.3.2. Underlying mechanisms

Figure 15 presents the Mantel test results for the four spatiotemporal congestion patterns across the two clusters. Green lines ($p \in [0.01, 0.05]$) and red lines ($p < 0.01$) represent significant relationship between the TTI values and the built environment features. The results support several key findings. First, on urban main arterials, road density, number of traffic lights, and distance to the nearest transit hubs are identified as primary congestion drivers. In contrast, on city expressways, ramp density and number
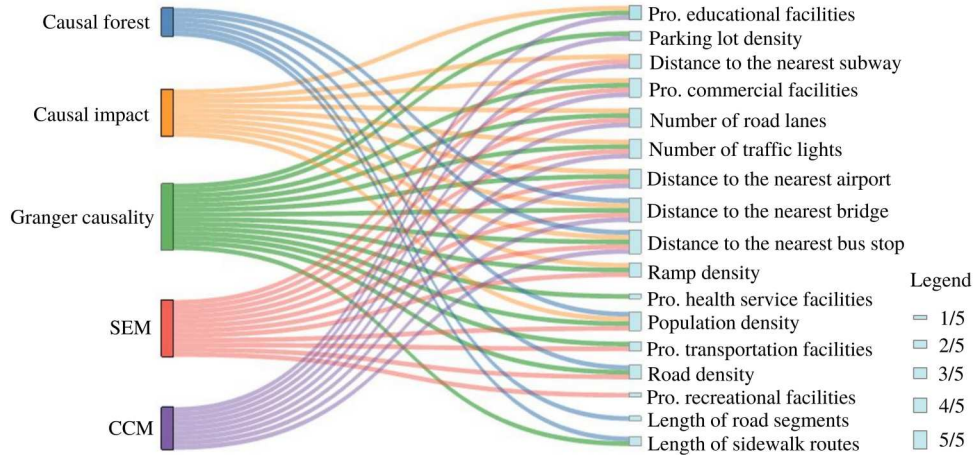
**Table 6.** Correlation and causality-based recommendations.

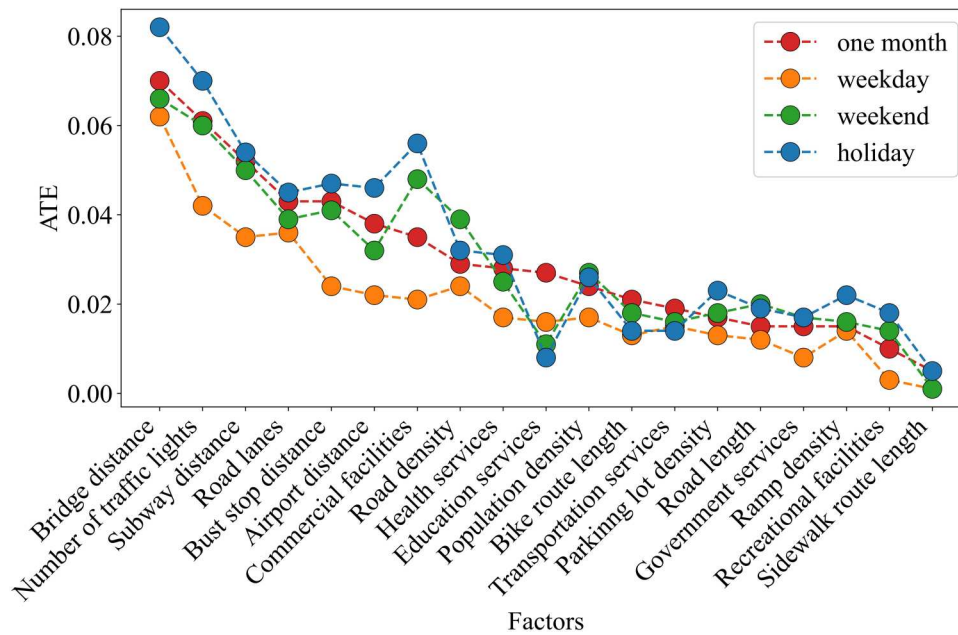| 'Five Ds' | Correlation-based suggestions and consequences | Causality-based findings and suggestions |
|---|---|---|
| 'Diversity' vs 'Density' | Increasing land use 'Diversity' in high 'Density' areas is thought to reduce trips to the city center (Song et al. 2019) – this can increase local congestion by complicating short-distance travel patterns. | 'Diversity' shows a positive causal effect, with time-dependent variations in its impact strength. This supports time-sensitive interventions, e.g. commute-priority corridors and commercial-flow lanes. |
| 'Density' vs 'Design' | Congestion caused by dense employment zones should be offset by expanding public transit (Bao et al. 2023) – this may induce new congestion hotspots near transit hubs. | 'Design' exerts a stronger causal effect than 'Density'. This support 'Design'-focused interventions, such as pressure-sensitive traffic-lights, and freight tunnel near bridge-related congestion hotspots. |
| 'Distance & Destination accessibility' | Public transit accessibility is viewed as critical for congestion mitigation (Ding et al. 2025) – increasing parking near public transit hubs may initially worsen congestion until key infrastructure thresholds are crossed. | This dimension presents is the most robust influence. It should be coupled with complementary measures, such as congestion pricing and expanded bike-sharing infrastructure. |

**Figure 11.** Parameter settings of the CCM.



**Figure 12.** Overall result. One factor is defined as a causal factor only if it is validated by at least three models. This leads to 11 causal factors.

of road lanes emerge as significant causal factors. Notably, distance to the nearest bridge appears as the most consistent key influential factor across both two roadway types, because it reflects network edges connecting distinct traffic communities (Sun et al. 2014). Second, non-weekday patterns show stronger causal associations with traffic lights and the proportion of public facilities, compared with weekday patterns. Additionally, the causal influence of first 'D' (Diversity) is found to be are weaker than that of the third 'D' (Design). Overall, the underlying causes of congestion on main urban arterials appear more complex than those on city expressways, potentially due to the complex geographical environment surrounding them.

**Figure 13.** Temporal sensitivity test of 19 correlated factors. The top-six factors exhibit temporal stability.

## 4. Discussion

### 4.1. Findings and recommendations

This study identifies causal explanations of urban congestion, thereby extending and refining those offered by previous correlation studies (Table 6). First, earlier correlational findings point to the need to increase land use 'Diversity' in high 'Density' areas (Song et al. 2019) and that 'Density'-induced congestion should be offset by public transit (Bao et al. 2023). Our study show that the causal effect of 'Design' is stronger than that of 'Diversity' and 'Density' in a mature city such as NYC. Second, the 'Distance & Destination accessibility' component shows consistently robust causal effects with Ding et al. (2025). Third, our identification of road density and traffic lights as having bidirectional causal relationships with congestion extends the unidirectional relationships reported by Duan et al. (2020).

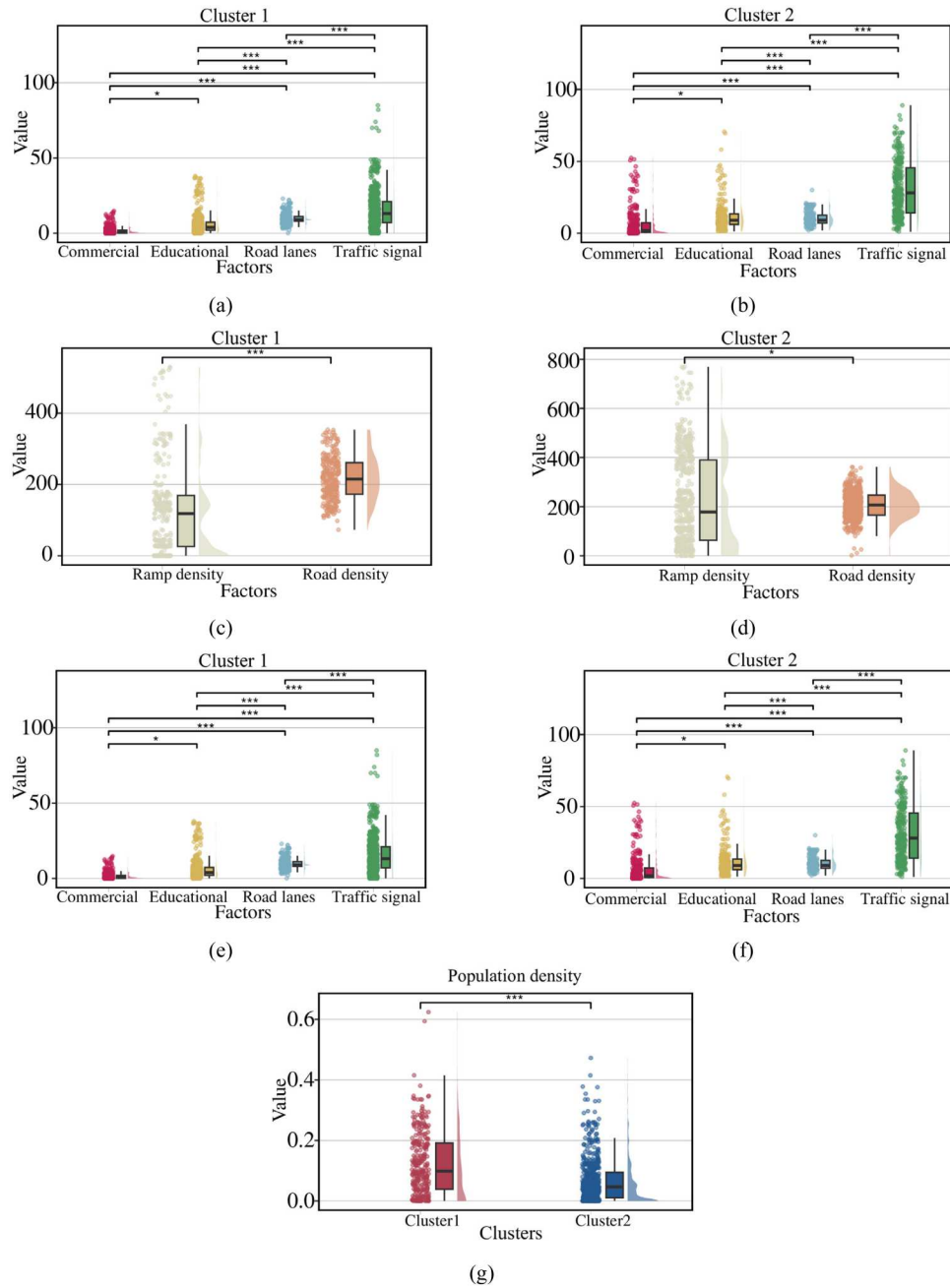### 4.2. Multi-model selection suggestions

This study integrates five complementary models to produce novel findings. However, these findings must be interpreted in conjunction with the specific application contexts suited to a multi-model approach, as summarized in Table 7. To ensure robustness, we recommend resolving conflicts through cross-validation: (i) verifying whether the data satisfy model assumptions, and (ii) conducting parameter sensitivity test.

### 4.3. Limitations and prospects

This study has three principal limitations that warrant discussion. First, excludes critical random variables that may influence traffic patterns and user behavior, such as real-time weather fluctuations and incident

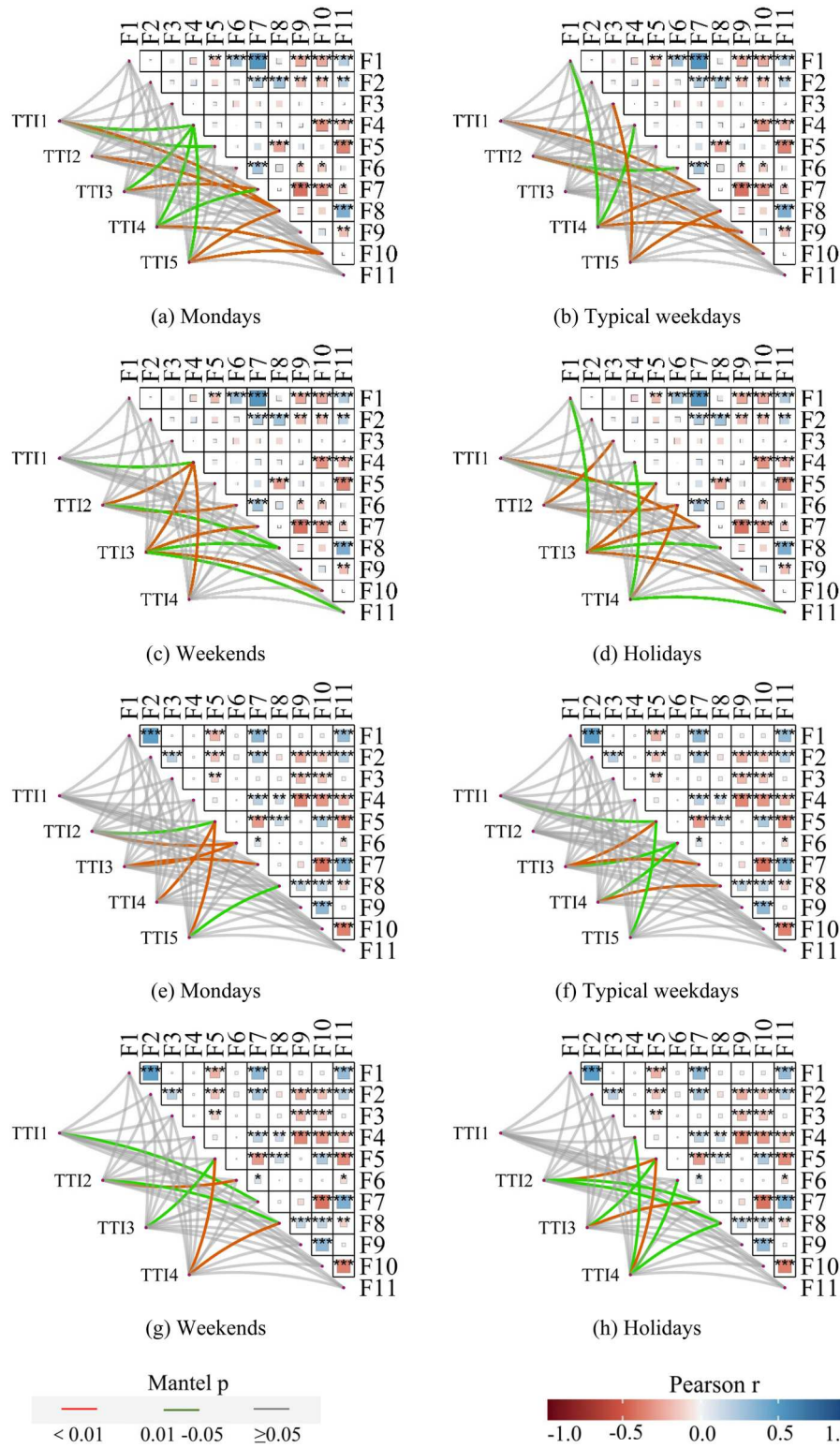**Table 7.** Contextualized multi-model strategy.

| Context | Recommended model combination | Advantages |
|---|---|---|
| Short-term assessment of transportation infrastructure | Granger causality + Causal Impact | Capture immediate causal effects, e.g. the short-term impact of newly built infrastructure. |
| Long-term land use planning | Causal Forest + CCM | Accounts for spatial heterogeneity and nonlinear, delayed effects. |
| Citywide congestion control | Granger causality + SEM | Identifies causal pathways, e.g. how bus stops affect congestion indirectly via transit hub usage. |
| Place-specific congestion control | Causal Forest + Causal Impact + SEM | Precisely detects spatial hotspots, while revealing causal pathways, e.g. the impact of bike lanes on congestion near schools. |

**Figure 14.** Statistical information of 11 factors in two types of perpetually congested roadways.

data. These observed factors were not controlled for potentially introducing heterogeneity into the estimation of interactions between the built environment and traffic congestion. Second, although we conducted sensitivity analysis comparing weekdays, weekends, and holidays, the exclusive reliance on data from December 2018 imposes temporal constraints. This limits the generalizability of our findings and precludes analysis of seasonal or interannual variation. Third, by focusing on the static 'five Ds' of the built environment factors, the current modeling framework does not account for the emerging sixth 'D'-the dynamic dimension, represented by demand responsive services. To address these limitations, future research should pursue three directions: (1) integrate or control for random exogenous factors; (2) utilize multi-seasonal and multi-annual data to capture temporal dynamics, including seasonal patterns and extreme weather effects; and (3) develop hybrid models that systematically incorporate dynamic the sixth 'D' components. These enhancements will improve the explanatory power and practical relevance of the framework.

**Figure 15.** Mantel test for four patterns of cluster 1 (a-d) and cluster2 (e-h). F1-F11 represent the 11 causal factors (Appendix Table 6). TTI1- TTI5 represents the average TTI at five periods (1:00-6:00, 7:00-9:00, 10:00-15:00, 16:00-19:00, 20:00-24:00).

## 5. Conclusion

The objective of this study is to investigate the causal relationship between traffic congestion and the built environment in New York City through three efforts: (1) estimating a multivariable least squares to identify

built environment factors significantly correlated with congestion; (2) establishing causal relationships between these factors using five causal inference models, while discussing a contextualized multi-model strategy for applying causal inference methods into urban research; and (3) examining the underlying impact mechanisms of four spatiotemporal patterns of perpetually congested roadways using Mantel test. The results highlight the limitations of relying solely on correlational analysis-namely, the risk of misleading feature selection and ineffective policy interventions. Based on the identified causal mechanisms, we propose targeted built environment interventions aimed at effectively mitigating traffic congestion.

## Through these efforts and results, our paper offers three contributions:

(1) The development of a framework for analyzing the causal importance of built environment factors, enabling the quantification of the global importance of the 'five Ds'. The New York City case study results ranked their significance as follows: Distance & Destination accessibility, Design, Diversity, and Density. Notably, public transit accessibility and traffic signal design emerge as critical for mitigating traffic congestion.

(2) The identification of distinct casual factors driving congestion on perpetually congested roadways-key bottlenecks in urban traffic networks-through temporal analysis across five daily time periods. Bridge-related congestion shows the strongest and most consistent impact.

(3) The provision of empirical grounding for several evidence-based urban planning recommendations. First, enhance accessibility to major transit hubs and increase traffic capacity at critical bridge connections. We suggest to deploy bicycle-sharing systems specifically tailored to first-and-last-mile connectivity at key transit nodes. Second, implement intelligent demand management vis congestion pricing at tunnel and bridge approaches to urban centers. Third, deploy adaptive traffic signal control systems that automatically trigger traffic signal extension when congestion exceeds defined thresholds. This approach can help break the vicious cycle of 'traffic flow aggregation, signal delay, and congestion exacerbation' (Naeem et al. 2024).

## Declaration of generative AI and AI-assisted technologies in writing

During the preparation of this work the authors used ChatGPT for editing, refining language, and enhancing clarity. After using this service, the authors reviewed and edited the content as needed and take full responsibility for the content of the published article. Note: Signifiance codes: *** $P<0.001$; **$P<0.01$; * $P<0.05$.

## Author contributions

W. Huan conceptualized the study, drafted the manuscript, and led the overall direction and planning of the paper. W. Huang conceptualized the study, contributed to the development of the theoretical framework and provided critical revisions of the manuscript. S. Li, X. Liu, and H. Wu participated in the interpretation of the results and offered significant insights during the revision process. M. Diao, H. Li, A. Grinberger, H. Liu, and C. Liu contributed significantly to the writing and editing process. All authors reviewed and approved the final version of the manuscript. Weihua Huan completed this work during a transfer from Tongji University (primary affiliation) to The Hong Kong Polytechnic University (secondary affiliation), and the research was conducted under the academic framework of both institutions. Special thanks are extended to the Department of Land Surveying and Geo-Informatics (LSGI) at The Hong Kong Polytechnic University for their academic and resource support.

## Disclosure statement

No potential conflict of interest was reported by the author(s).

## Funding

## Data availability statement

The publicly accessible average hourly travel speed data were downloaded from Uber Movement (https://movement.uber.com/explore/new_york/speeds/query?dt[tpb]=ALL_DAY&dt[wd;]=1,2,3,4,5,6,7&dt[dr][sd]=2018-11-28&dt[dr][ed]=2018-12-30&ff=&lat.=40.7491664&lng.=-74.0154715&z.=11.9&lang=en-US.) on April 20, 2023. As of October 1, 2023, Uber Movement had ceased providing services. Data used in this study are available upon request. The publicly accessible built environment data were obtained from NYC Open Data (https://opendata.cityofnewyork.us/), using the following specific datasets:

(1) Points of Interest: https://data.cityofnewyork.us/City-Government/Points-Of-Interest/rxuy-2muj;
(2) Building Footprints: https://data.cityofnewyork.us/Housing-Development/Building-Footprints/nqwf-w8eh;
(3) Bus Stop Shelters: https://data.cityofnewyork.us/Transportation/Bus-Stop-Shelters/qafz-7myz;
(4) Subway Stations: https://data.cityofnewyork.us/Transportation/Subway-Stations/arq3-7z49;
(5) Parking Lots: https://data.cityofnewyork.us/City-Government/Parking-Lot/h7zy-iq3d;
(6) Leading Pedestrian Interval Signals: https://data.cityofnewyork.us/Transportation/VZV_Leading-Pedestrian-Interval-Signals/mqt5-ctec;
(7) Bike Routes: https://data.cityofnewyork.us/Transportation/New-York-City-Bike-Routes/7vsa-caz7;
(8) Truck Routes: https://data.cityofnewyork.us/Transportation/New-York-City-Truck-Routes-Map-/wnu3-egq7;
(9) Sidewalks: https://data.cityofnewyork.us/City-Government/Sidewalk/vfx9-tbb6;
(10) Bridge Ratings: https://data.cityofnewyork.us/Transportation/Bridge-Ratings/4yue-vjfc;
(11) Airport Polygons: https://data.cityofnewyork.us/City-Government/Airport-Polygon/xfhz-rhsk.

The datasets listed above may have been updated over time; however, the average travel speed data and built environment data used in this study were temporally aligned to ensure consistency.

## ORCID

*Xintao Liu*   http://orcid.org/0000-0002-7323-9878

## References

Albini, E., J. Long, D. Dervovic, and D. Magazzeni. 2022. "Counterfactual Shapley Additive Explanations." Proceedings of the 2022 ACM Conference on Fairness, Accountability, And Transparency (pp. 1054-1070). https://doi.org/10.1145/3531146.353316.

Bao, Z. K., S. T. Ng, G. Yu, X. L. Zhang, and Y. F. Ou. 2023. "The Effect of the Built Environment on Spatial-Temporal Pattern of Traffic Congestion in a Satellite City in Emerging Economies." *Developments in the Built Environment* 14:100173. https://doi.org/10.1016/j.dibe.2023.100173.

Bao, Z. K., Y. Y. Ou, S. Z. Chen, and T. Wang. 2022. "Land use Impacts on Traffic Congestion Patterns: A Tale of a Northwestern Chinese City." *Land* 11 (12): 2295. https://doi.org/10.3390/land11122295.

Benito-Moreno, M., J. Carpio-Pinedo, and P. J. Lamíquiz-Daudén. 2025. "Proximity Features: A Random Forest Approach to the Influence of the Built Environment on Local Travel Behavior." *Urban Science* 9 (4): 122. https://doi.org/10.3390/urbansci9040122.

Brodersen, K. H., F. Gallusser, J. Koehler, N. Remy, and S. L. Scott. 2015. "Inferring Causal Impact Using Bayesian Structural Time-series Models." The Annals of Applied Statistics 9 (1): 247–274. https://doi.org/10.1214/14-AOAS788.

Cervero, R., and K. Kockelman. 1997. "Travel Demand and the 3Ds: Density, Diversity, and Design." *Transportation Research Part D: Transport and Environment* 2 (3): 199–219. https://doi.org/10.1016/S1361-9209(97)00009-6.

Cook, C., A. Kreidieh, S. Vasserman, H. Allcott, N. Arora, F. van Sambeek, A. Tomkins, and E. Turkel. 2025. *The Short-Run Effects of Congestion Pricing in New York City (No. w33584)*, Technical report. Massachusetts: National Bureau of Economic Research. https://doi.org/10.3386/w33584.

Deng, X., J. Zhang, S. Liao, C. Zhong, F. Gao, and L. Teng. 2022. "Interactive Impacts of Built Environment Factors on Metro Ridership Using GeoDetector: From the Perspective of TOD." *ISPRS International Journal of Geo-Information* 11 (12): 623. https://doi.org/10.3390/ijgi11120623.

Diao, M., H. Kong, and J. H. Zhao. 2021. "Impacts of Transportation Network Companies on Urban Mobility." *Nature Sustainability* 4 (6): 494–500.

Ding, H., Z. Zhao, S. Wang, Y. Zhang, X. Zheng, and X. Lu. 2025. "Quantifying the Impact of Built Environment on Traffic Congestion: A Nonlinear Analysis and Optimization Strategy for Sustainable Urban Planning." *Sustainable Cities and Society* 122:106249. https://doi.org/10.1016/j.scs.2025.106249.

Dormann, C. F., J. M. Calabrese, G. Guillera-Arroita, E. Matechou, V. Bahn, K. Bartoń, C. M. Beale, et al. 2018. "Model Averaging in Ecology: A Review of Bayesian, Information-Theoretic, and Tactical Approaches for Predictive Inference." *Ecological Monographs* 88 (4): 485–504. https://doi.org/10.1002/ecm.1309.

Duan, X., J. Xu, Y. Chen, and R. Jiang. 2020. "Analysis of Influencing Factors on urban Traffic Congestion and Prediction of Congestion Time Based on Spatiotemporal Big Data." 2020 International Conference on Big Data, Artificial Intelligence and Internet of Things Engineering (ICBAIE) (pp. 75-78). IEEE. https://doi.org/10.1109/ICBAIE49996.2020.00022.

Ewing, R., and R. Cervero. 2010. "Travel and the Built Environment: A Meta-analysis." *Journal of the American Planning Association* 76 (3): 265–294. https://doi.org/10.1080/01944361003766766.

Gao, C., X. Lai, S. Li, Z. Cui, and Z. Long. 2023. "Bibliometric Insights into the Implications of Urban Built Environment on Travel Behavior." *ISPRS International Journal of Geo-Information* 12 (11): 453. https://doi.org/10.3390/ijgi12110453.

Golob, T. F. 2003. "Structural Equation Modeling for Travel Behavior Research." *Transportation Research Part B: Methodological* 37 (1): 1–25. https://doi.org/10.1016/S0191-2615(01)00046-7.

Huang, G. X., and D. Xu. 2023. "The Last Mile Matters: Impact of Dockless Bike-Sharing Services on Traffic Congestion." *Transportation Research Part D: Transport and Environment* 121:103836. https://doi.org/10.1016/j.trd.2023.103836.

Imbens, G. W., and D. B. Rubin. 2015. *Causal Inference in Statistics, Social, and Biomedical Sciences.* Cambridge: Cambridge University Press.

Jiang, F., L. Ma, T. Broyd, W. Chen, and H. Luo. 2022. "Digital Twin Enabled Sustainable Urban Road Planning." *Sustainable Cities and Society* 78:103645. https://doi.org/10.1016/j.scs.2021.103645.

Ju, H., Z. Zhang, L. Zuo, J. Wang, S. Zhang, X. Wang, and X. Zhao. 2016. "Driving Forces and Their Interactions of Built-up Land Expansion Based on the Geographical Detector–a Case Study of Beijing, China." *International Journal of Geographical Information Science* 30 (11): 2188–2207. https://doi.org/10.1080/13658816.2016.1165228.

Kaiser, Z. A., and A. Deb. 2025. "Sustainable Smart City and Sustainable Development Goals (SDGs): A Review." *Regional Sustainability* 6 (1): 100193.

Kamat, P. V. 2025. "Correlation, Causation and Comparison." *ACS Energy Letters* 10 (3): 1540–1541. https://doi.org/10.1021/acsenergylett.5c00631.

Kong, X., J. Yang, and Z. Yang. 2015. "Measuring Traffic Congestion with Taxi GPS Data and Travel Time Index." *CICTP*: 3751–3762. https://doi.org/10.1061/9780784479292.346.

Koźlak, A., and D. Wach. 2018. "Causes of Traffic Congestion in Urban Areas. Case of Poland." *SHS Web of Conferences. EDP Sciences* 57:01019. https://doi.org/10.1051/shsconf/20185701019.

Lee, S., S. Lee, and D. W. Putri. 2025. "Multifaceted Associations between Built Environments and POI Visit Patterns by Trip Purposes." *Cities* 161: 105903. https://doi.org/10.1016/j.cities.2025.105903.

Li, K. T. 2020. "Statistical Inference for Average Treatment Effects Estimated by Synthetic Control Methods." *Journal of the American Statistical Association* 115 (532): 2068–2083. https://doi.org/10.1080/01621459.2019.1686986.

Li, Z., C. Liang, Y. L. Hong, and Z. J. Zhang. 2022. "How Do on-Demand Ridesharing Services Affect Traffic Congestion? The Moderating Role of Urban Compactness." *Production and Operations Management* 31 (1): 239–258. https://doi.org/10.1111/poms.13530.

Li, D., J. Ma, T. Cheng, J. L. van Genderen, and Z. Shao. 2019. "Challenges and Opportunities for the Development of MEGACITIES." *International Journal of Digital Earth* 12 (12): 1382–1395. https://doi.org/10.1080/17538947.2018.1512662.

Li, C., D. Wang, H. Chen, and E. Liu. 2024. "Analysis of Urban Congestion Traceability: The Role of the Built Environment." *Land* 13 (2): 255. https://doi.org/10.3390/land13020255.

Liang, Y., B. J. Yu, X. J. Zhang, Y. Lu, and L. C. Yang. 2023. "The Short-Term Impact of Congestion Taxes on Ridesourcing Demand and Traffic Congestion: Evidence from Chicago." *Transportation Research Part A: Policy and Practice* 172:103661. https://doi.org/10.1016/j.tra.2023.103661.

Liu, J., J. Cui, L. Xiao, D. Lin, and L. Yang. 2025. "Non-linear Built Environment Effects on Travel Behavior Resilience under Extreme Weather Events." *Transportation Research Part D: Transport and Environment* 143:104753. https://doi.org/10.1016/j.trd.2025.104753.

Liu, F., and Y. Deng. 2020. "Determine the Number of Unknown Targets in Open World Based on Elbow Method." *Transactions on Fuzzy Systems* 29 (5): 986–995. https://doi.org/10.1109/TFUZZ.2020.2966182.

Managing urban traffic congestion. 2007. "Transport Research Centre, European Conference of Ministers of Transport." Available at: https://www.itf-oecd.org/sites/default/files/docs/07congestion.pdf.

Naeem, R., R. M. Shoaib, Z. Ahmed, and A. Razaq. 2024. "A Literature Survey on Signalized and non-signalized Corridors." *Jurnal Syntax Transformation* 5 (10): 1179–1187.

Olayode, I. O., A. Severino, F. J. Alex, and E. Jamei. 2025. "Traffic Flow Modelling of Vehicles on a Six Lane Freeway: Comparative Analysis of Improved Group Method of Data Handling and Artificial Neural Network Model." *Results in Engineering* 25: 104094. https://doi.org/10.1016/j.rineng.2025.104094.

Oti, E., and M. Olusola. 2024. ""Overview of Agglomerative Hierarchical Clustering Methods." British Journal of Computer." *Networking and Information Technology* 7:14–23. https://doi.org/10.52589/BJCNIT-CV9POOGW.

Pan, Y., S. Chen, S. Niu, Y. Ma, and K. Tang. 2020. "Investigating the Impacts of Built Environment on Traffic States Incorporating Spatial Heterogeneity." *Journal of Transport Geography* 83: 102663. https://doi.org/10.1016/j.jtrangeo.2020.102663.

Pi, M., H. Yeon, H. Son, and Y. Jang. 2019. "Visual Cause Analytics for Traffic Congestion." *IEEE Transactions on Visualization and Computer Graphics* 27 (3): 2186–2201. https://doi.org/10.1109/TVCG.2019.2940580.

Qian, X., T. Lei, J. Xue, Z. Lei, and S. V. Ukkusuri. 2020. "Impact of Transportation Network Companies on Urban Congestion: Evidence from Large-Scale Trajectory Data." *Sustainable Cities and Society* 55:102053.

Rahman, M. M., P. Najaf, M. G. Fields, and J. C. Thill. 2022. "Traffic Congestion and Its Urban Scale Factors: Empirical Evidence from American Urban Areas." *International Journal of Sustainable Transportation* 16 (5): 406–421. https://doi.org/10.1016/j.scs.2020.102053.

Rotari, M., and M. Kulahci. 2024. "Correlation to Causality. Quality Engineering." *Quality Engineering* 37 (1): 162–172. https://doi.org/10.1080/08982112.2024.2372489.

Rousseeuw, P. J. 1987. "Silhouettes: A Graphical aid to the Interpretation and Validation of Cluster Analysis." *Journal of Computational and Applied Mathematics* 20:53–65. https://doi.org/10.1016/0377-0427(87)90125-7.

Saarela, M. 2024. "On the Relation of Causality-Versus Correlation-Based Feature Selection on Model Fairness." In Proceedings of the 39th ACM/SIGAPP Symposium on Applied Computing pp. 56-64. https://doi.org/10.1145/3605098.363601.

Saberi, M., H. Hamedmoghadam, M. Ashfaq, S. A. Hosseini, Z. Gu, S. Shafiei, D. J. Nair, et al. 2020. "A Simple Contagion Process Describes Spreading of Traffic Jams in Urban Networks." *Nature Communication* 11 (1): 1616.

Shen, T., Y. Hong, M. M. Thompson, J. P. Liu, and X. P. Huo. 2020. "How Does Parking Availability Interplay with the Land use and Affect Traffic Congestion in Urban Areas? The Case Study of Xi'an, China." *Sustainable Cities and Society* 57:102126. https://doi.org/10.1016/j.scs.2020.102126.

Shojaie, A., and E. B. Fox. 2022. "Granger Causality: A Review and Recent Advances." *Annual Review of Statistics and Its Application* 9 (1): 289–319. https://doi.org/10.1146/annurev-statistics-040120-010930.

Shrestha, N. 2020. "Detecting Multicollinearity in Regression Analysis." *American Journal of Applied Mathematics and Statistics* 8 (2): 39–42. https://doi.org/10.12691/ajams-8-2-1.

Solecki, W., and C. Rosenzweig. 2019. "New York city Panel on Climate Change 2019 Report Chapter 9: Perspectives on a City in a Changing Climate 2008-2018." Annals of the New York Academy of Sciences, 1439(GSFC-E-DAA-TN66874). https://doi.org/10.1111/nyas.14017.

Somers, K. M., and D. A. Jackson. 2022. "Putting the Mantel Test Back Together Again." *Ecology* 103 (10): e3780. https://doi.org/10.1002/ecy.3780.

Song, J. C., C. L. Zhao, S. P. Zhong, and T. A. S. Nielsen. 2019. ""Mapping Spatio-Temporal Patterns and Detecting the Factors of Traffic Congestion with Multi-source Data Fusion and Mining Techniques." Computers." *Environment and Urban Systems* 77:101364. https://doi.org/10.1016/j.compenvurbsys.2019.101364.

Su, P., Y. Yan, H. Li, H. Wu, C. Liu, and W. Huang. 2025. "Images and Deep Learning in Human and Urban Infrastructure Interactions Pertinent to Sustainable Urban Studies: Review and Perspective." *International Journal of Applied Earth Observation and Geoinformation* 136: 841–854. https://doi.org/10.1016/j.jag.2024.104352.

Sun, H., J. Wu, D. Ma, and J. Long. 2014. "Spatial Distribution Complexities of Traffic Congestion and Bottlenecks in Different Network Topologies." *Applied Mathematical Modelling* 38 (2): 496–505. https://doi.org/10.1016/j.apm.2013.06.027.

Tay, R. 2017. "Correlation, Variance Inflation and Multicollinearity in Regression Model." *Journal of the Eastern Asia Society for Transportation Studies* 12:2006–2015. https://doi.org/10.11175/easts.12.2006.

Tracy, A. J., P. Su, A. W. Sadek, and Q. Wang. 2011. "Assessing the Impact of the Built Environment on Travel Behavior: A Case Study of Buffalo, New York." *Transportation* 38:663–678. https://doi.org/10.1007/s11116-011-9337-x.

Tsonis, A. A., E. R. Deyle, H. Ye, and G. Sugihara. 2018. "Convergent Cross Mapping: Theory and an Example." In *Advances in Nonlinear Geosciences*, edited by A. A. Tsonis, 587–600. Cham: Springer. https://doi.org/10.1007/978-3-319-58895-7_27.

Wager, S., and S. Athey. 2018. "Estimation and Inference of Heterogeneous Treatment Effects Using Random Forests." *Journal of the American Statistical Association* 113 (523): 1228–1242. https://doi.org/10.1080/01621459.2017.1319839.

Wang, D. G., and M. Zhou. 2017. "The Built Environment and Travel Behavior in Urban China: A Literature Review." *Transportation Research Part D: Transport and Environment* 52:574–585. https://doi.org/10.1016/j.trd.2016.10.031.

Wiener, N. 1956. "What Is Information Theory." *Transactions on Information Theory* 2 (2): 48.

Zhang, K. S., D. Sun, S. W. Shen, and Y. Zhu. 2017a. "Analyzing Spatiotemporal Congestion Pattern on Urban Roads Based on Taxi GPS Data." *Journal of Transport and Land Use* 10 (1): 675–694.

Zhang, T. Q., L. S. Sun, L. Y. Yao, and J. Rong. 2017b. "Impact Analysis of Land use on Traffic Congestion Using Real-Time Traffic and POI." *Journal of Advanced Transportation* 2017 (1): 7164790. https://doi.org/10.1155/2017/7164790.

Zhang, W., B. Sun, and C. Zegras. 2021. "Sustainable Built Environment and Travel Behavior: New Perspectives, new Data, and new Methods." *Transportation Research Part D: Transport and Environment* 97:102966. https://doi.org/10.1016/j.trd.2021.102966.

Zhu, H., H. Guan, Y. Han, and W. Li. 2019. "A Study of Tourists' Holiday Rush-Hour Avoidance Travel Behavior Considering Psychographic Segmentation." *Sustainability* 11 (13): 3755. 10.3390/su11133755.