
Explainable Spatio-Temporal Forecasting with Shape Functions

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Spatio-temporal modeling and forecasting are challenging due to their complicated
2 spatial dependence, temporal dynamics, and scenarios. Many statistical models,
3 such as Spatial Auto-regression Model (SAR) and Spatial Dynamic Panel Data
4 Model (SDPD), are restricted by a pre-specified spatial weight matrix and thus
5 are limited to reflect its flexibility. Graph-based or convolution-based methods
6 can learn more flexible representations, but they fail to show the exact interactions
7 between locations due to the lack of explainability. This paper proposes a spatial re-
8 gression model with shape functions to address the limitations of existing methods.
9 Our method learns the shape functions by incorporating shape constraints, which
10 are able to capture spatial variability or distance-based effects over distance. There-
11 fore, our approach enjoys a learnable spatial weight matrix with a distance-based
12 explanation. We demonstrate our method’s efficiency and forecasting performance
13 on synthetic and real data.

14 1 Introduction

15 Spatio-temporal data is widely observed in many areas, such as transportation (33; 27), climatology
16 (2), and environmental research(19). The popularity of spatio-temporal data brings varieties of tasks
17 for researchers, and one of the key tasks is forecasting. Spatio-temporal data has some inherent
18 characteristics, namely, spatial dependence and temporal dynamics, which need to be considered for
19 modeling and forecasting.

20 Spatial dependence means that the observations at different locations are not independent, and
21 observations at closer locations often have a stronger correlation. In the statistics community,
22 extensive research has been conducted to model spatial dependence, and various spatial models have
23 been proposed. For example, in the spatial autoregressive (SAR) models, the spatial dependence is
24 modeled by a product of an unknown parameter and a pre-specified spatial weight matrix (4; 1; 11; 12).
25 Combined with the panel data, various types of spatial panel data models have been used to analyze
26 spatio-temporal data (35; 13; 7; 22). One limitation of the autoregressive models is that the elements
27 of the spatial weight matrix are pre-specified, such as an inverse distance. Although these pre-
28 specified spatial weight matrices are applied to capture decreased distance-based effects, they fail to
29 capture complex distance relations in real-world applications.

30 Researchers in the computer science community have developed various methods modeling spatio-
31 temporal data using deep neural networks. Various neural network architectures have been proposed
32 and applied to spatio-temporal forecasting, for example, spatio-temporal LSTM (31), fully connected
33 gated graph architecture (20), Convolutional LSTM (23) and etc. One advantage of these methods
34 is that they can incorporate unstructured data and rely on a high-performance computing platform
35 to learn complicated representations for spatio-temporal problems. However, a critical limitation of
36 these methods is that they fail to explain how the spatial interaction works explicitly. The lack of

37 interpretability restricts its reliability and deep insights into the underlying spatio-temporal process.
38 The explanation can be obtained if we can estimate the coefficient matrix that intuitively explains
39 spatio-temporal interactions.

40 In this paper, we propose an Explainable Spatio-Temporal Forecasting (ESTF) model, which utilizes
41 a spatial autoregressive model with shape functions to address the current limitations. Our method
42 extends the vector autoregressive (VAR) model (24) by incorporating distance information into the
43 temporal coefficient matrix using shape functions (3). The shape constraints are designed to be
44 consistent with the common fact that observations from neighbours have stronger spatial dependence
45 versus long-distance pairs. It is known as Tobler's First Law, which is "Everything is related to
46 everything else, but near things are more related than distant things"(26; 18). Unlike the pre-specified
47 spatial weight matrix, this coefficient matrix is learnable and is thus more flexible in capturing
48 real-world complex spatial relations. Moreover, the shape functions are represented as a combination
49 of basis functions, and thus a smaller number of parameters needs to be estimated. Finally, ESTF can
50 be easily extended to forecasting in non-stationary scenarios using a dynamic spatial weight matrix.
51 We conduct experiments on both simulated and real data, and the results demonstrate that our method
52 achieves better forecast accuracy and is computationally efficient and more explainable.

53 2 Related work

54 **Statistical models** Several works focus on temporal dynamics when considering spatio-temporal
55 forecasting problems. The classical time series models, such as VAR, and ARIMA models, are applied
56 to spatio-temporal process modeling(21; 38). Besides, a spatial weight matrix is also introduced to the
57 ARIMA model to capture spatial dependence (28). The non-stationarity, particularly unit-root non-
58 stationarity, is mainly modeled by ARIMA or Co-integration models. In addition, spatial regression
59 models or panel data are classical models in econometrics and can also be applied to model spatio-
60 temporal problems. These models, for example, spatial auto-regression models, take spatial weight
61 matrix into consideration and estimate parameters in the framework of regression. However, the
62 common characteristics of these models need a pre-specified spatial weight matrix(35; 6). Elements
63 in the matrices are generally an inverse distance of corresponding locations. Meanwhile, these
64 spatial models focus on statistical inference on the scalar parameters placed before the spatial weight
65 matrix(25). Although there are many choices for the spatial weight matrix, such as inverse distance,
66 adjacency relationships, and K-nearest neighbors, there is a lack of research on estimating the spatial
67 weight matrix. The pre-specified spatial weight matrix restricts models' application and fails to
68 capture more complicated underlying spatial dependence. Some researchers developed a sparse
69 spatio-temporal model that can estimate a sparse spatial weight matrix (17). The strict sparse setting
70 also restricts the wide application of the spatial weight matrix.

71 **Graph-based methods** Graph-based methods are widely applied for a non-Euclidean domain.
72 Some types of spatio-temporal data, for example, traffic flow data or brain network data, can be
73 represented as graphs. The graph structures well model the complicated spatial dependence. Thus,
74 the definition or pre-specified graphs structure is normally required when developing a graph-based
75 model. Related works can be found in (30; 14). The common typical method is GraphCNN, which is
76 to apply a convolutional transformation to the neighbors of each node (29; 34). The graph convolution
77 can capture patterns and features in the spatial domain. Graph-based methods have been proposed
78 and widely applied to lots of real cases. Traffic flow data modeling and forecasting is a popular topic
79 in this area (30; 20). Other topics, for example, climate sensor data (16), video (10) and etc, are also
80 applied by variant graph-based models. RNN or LSTM combined with graphs, i.e., a sequence of
81 graphs, are also considered in spatio-temporal forecasting problems (10).

82 **CNN-based methods** Unlike graph-based methods, CNN-based methods are more suitable for
83 modeling spatio-temporal data collected in regular grid locations. It applies filters to find relationships
84 between neighboring inputs. Although some works (32) applied convolution neural networks to
85 model non-grid traffic data, it is more common to see CNN-based methods process grid structures,
86 e.g., images, video rather than a general domain. As some spatio-temporal data are collected from a
87 regular grid in the Euclidean space (29), they thus can be viewed as a kind of special image. The CNN
88 structure combined with RNN or LSTM has been developed to make forecasting for spatio-temporal
89 data, for example, diffusion convolutional RNN (15), Convolutional LSTM networks (23; 36)and etc.

90 3 Proposed method

91 3.1 Problem formulation and notation

92 We use a $n \times 1$ vector $\mathbf{X}_t = \{\mathbf{x}_{1t}, \mathbf{x}_{2t}, \dots, \mathbf{x}_{nt}\}$ to denote observations at time t , where n is
 93 the number of locations. At each location i , $\mathbf{S}_i = (\mathbf{c}_i^x, \mathbf{c}_i^y)$ is the coordinates of the location i .
 94 The distance between location \mathbf{S}_i and \mathbf{S}_j is $d_{ij} = \sqrt{(d_{ij}^x)^2 + (d_{ij}^y)^2}$, where $d_{ij}^x = |c_i^x - c_j^x|$ and
 95 $d_{ij}^y = |c_i^y - c_j^y|$. Our goal is to make forecasting for spatio-temporal data: given training data set
 96 $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_T$, we would like to make forecasting for the next h , $\hat{\mathbf{X}}_{T+1}, \dots, \hat{\mathbf{X}}_{T+h}$.

97 3.2 The stationary spatio-temporal model with shape functions

98 We first consider the stationary case. To model the spatio-temporal stationary process, we consider
 99 the following model

$$\mathbf{X}_t = \sum_{k=1}^P \mathbf{W}_k \mathbf{X}_{t-k} + \epsilon_t, \quad (1)$$

100 where \mathbf{W}_k is a spatial weight matrix for capturing the spatial dependence at lag k , and ϵ_t is white
 101 noise. Moreover, we assume the (i, j) th element of \mathbf{W}_k , $w_{ij}^{(k)}$, depends on the distance d_{ij} . That is,
 102 $w_{ij}^{(k)}$ depends on a function $f_k(d_{ij})$.

103 For spatio-temporal data, the spatial dependence, represented by $w_{ij}^{(k)}$, between locations decreases
 104 as the distance between two locations increases. In other words, there is a shape constraint for
 105 the function $f_k(d)$, such as a decreasing function. In order to estimate the shape function, we
 106 model $f_k(d)$ as a linear combination of basis functions $g_i(d)$, $i = 1, 2, \dots, m$. More specifically,
 107 the shape function $f_k(d)$ is a linear combination of basis functions and coefficients with positive
 108 value $f_k(d) = a_{1,k}^2 g_1(d) + \dots + a_{m,k}^2 g_m(d)$, where $a_{1,k}, \dots, a_{m,k}$ are parameters to be estimated.
 109 The constraint of decrease needs parameters non-negative and thus each parameters squared. The
 110 spatial weight matrix can take the value of decreased shape function directly. The element of \mathbf{W}_k
 111 is $w_{ij}^{(k)} = f_k(d_{ij})$. The details of the shape function and the corresponding basis functions can be
 112 found in Section 3.4

The parameters in shape functions can be estimated from the neural network illustrated in Figure 1.
 The neural network can be trained from the following criterion:

$$\min_{\{W_k\}_{k=1}^P} \sum_{t=1}^T \|\mathbf{X}_t - \hat{\mathbf{X}}_t\|^2 = \sum_{t=1}^T \|\mathbf{X}_t - \sum_{k=1}^P \hat{\mathbf{W}}_k \hat{\mathbf{X}}_{t-k}\|^2.$$

113 3.3 The non-stationary spatio-temporal model with time-variant shape functions

114 The static spatial weight matrix \mathbf{W}_k can reflect spatial dependence and thus can be applied to
 115 stationary scenarios. Next, we consider the nonstationary case. Therefore, we extend the stationary
 116 model to non-stationary cases. The spatial weight matrices only reflect static relationships across time
 117 lags in the static model. Unlike these settings, we change spatial weight matrices to be time-variant.
 118 The spatial weight matrices formed by time-variant shape functions can thus capture non-stationary
 119 dynamic spatial dependence. The non-stationary model has the form below,

$$\mathbf{X}_t = \sum_{k=1}^P \mathbf{W}_{t,k} \mathbf{X}_{t-k} + \epsilon_t. \quad (2)$$

where ϵ_t is white noise, and $\mathbf{W}_{t,k}$ relies on shape function $f_{t,k}(d)$. Similar with stationary settings,
 the time-variant shape functions are still represented as a linear combination of basis functions
 $g_i(d)$, $i = 1, 2, \dots, m$. The coefficients are therefore time-variant. The shape function at time t has
 the form below $f_{t,k}(d) = a_{1,t,k}^2 g_1(d) + \dots + a_{m,t,k}^2 g_m(d)$. Unlike stationary setting, the coefficients
 of nonstationary setting, $\{a_{i,t,k}\}_{i=1}^m$, depend on the time t . The non-stationary model can be trained
 from the criterion by minimizing

$$\min_{\{W_{t,k}\}_{k=1}^P} \|\mathbf{X}_t - \hat{\mathbf{X}}_t\|^2 = \|\mathbf{X}_t - \sum_{k=1}^P \hat{\mathbf{W}}_{t,k} \hat{\mathbf{X}}_{t-k}\|^2.$$

120 The networks for the stationary model as well as the non-stationary model are presented in the Figure
 121 1.

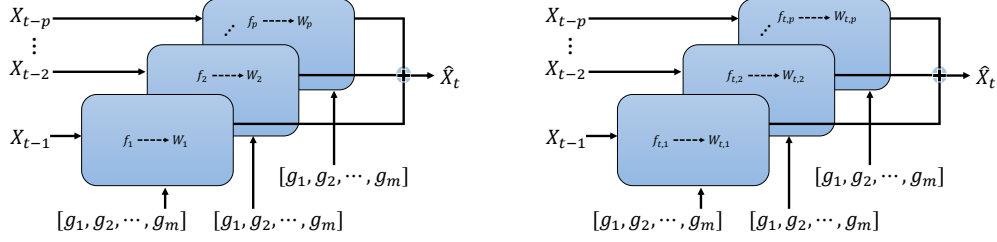


Figure 1: The neural network for the stationary spatio-temporal process (left) and non-stationary spatio-temporal process (right).

122 3.4 The basis functions for shape functions

123 The shape functions are integrated into our model to obtain distance-based explanations in stationary
 124 and non-stationary scenarios. The motivation of the proposed shape functions is that as the distance
 125 between two observations increases, the effects between these two locations decreases. These distance-
 126 based effects can be reflected in spatial weight matrix \mathbf{W} and each element in the matrix can measure
 127 how the corresponding locations interact. The shape function is represented as a linear combination of
 128 basis functions. The basis functions, satisfying shape constraint, rely on the corresponding definition
 129 of basis functions.

130 **Definition of basis functions for various shape constraints.** We list the definition of basis functions
 131 for increased and decreased shape (3). The distance quantile among $\{d_{i_1, j_1}, d_{i_2, j_2}, \dots, d_{i_N, j_N}\}$
 132 at quantile level q_1, q_2, \dots, q_m is denoted by $\{d_{(1)}, d_{(2)}, \dots, d_{(m)}\}$, where $0 \leq q_1 < q_2 < \dots < q_m \leq$
 133 1 and $\{q_1, q_2, \dots, q_m\} = \{\frac{1}{m}, \frac{2}{m}, \dots, 1\}$. Here, we can set the number of $m \ll n^2$, and thus, the
 134 number of parameters is significantly reduced.

135 For the constraint of monotone decreasing function, the basis function is defined as $g_i(d) = \mathbf{1}_{\{d < d_{(i)}\}}$.
 136 The basis function for the shape function with the constraint of concave decrease is defined as
 137 $g_i(d) = (d_{(i)} - d)\mathbf{1}_{\{d_{(i)} \leq d\}}$ and convex decrease is defined as $g_i(d) = (d_{(i)} - d)\mathbf{1}_{\{d \leq d_{(i)}\}}$, for
 138 $1 \leq i \leq m$. Figure 2 shows the definition of basis functions for monotone decreased and increased
 139 shape functions, respectively. We only present four basis functions for each shape and each of them
 140 is related to four quantile levels. The dashed lines indicate the turning points for each basis function
 141 and they equal one or zero at the beginning and turn to zero or one at turning points.

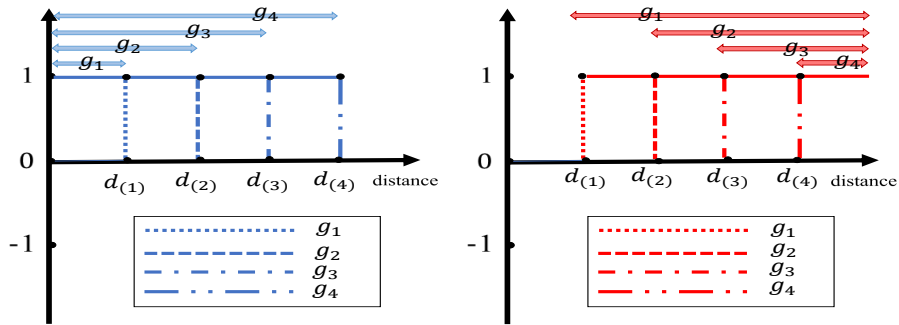


Figure 2: The basis functions for decreased shape (left) and for increased shape (right). The arrows indicate domain of each basis functions.

142 3.5 Model forecasting

143 The stationary model requires fixed shape functions and related spatial weight matrix are time-
 144 invariant. Given training data set $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_T$, we can estimate spatial weight matrix
 145 $\hat{W}_1, \hat{W}_2, \dots, \hat{W}_p$ and make forecasting iteratively. That is $\hat{\mathbf{X}}_{T+1} = \sum_{k=1}^p \hat{W}_k \mathbf{X}_{T+1-k}$, $\hat{\mathbf{X}}_{T+2} =$
 146 $\hat{W}_1 \hat{\mathbf{X}}_{T+1} + \sum_{k=2}^p \hat{W}_k \mathbf{X}_{T+2-k}, \dots, \hat{\mathbf{X}}_{T+h} = \sum_{k=1}^p \hat{W}_k \hat{\mathbf{X}}_{T+h-k}$.

The non-stationary model incorporate time-variant spatial weight matrix $\hat{W}_{t,\cdot}$. Given the training
 data set $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_T$, we can obtain corresponding shape functions $\hat{f}_{1,\cdot}, \hat{f}_{2,\cdot}, \dots, \hat{f}_{T,\cdot}$, where
 \cdot denotes time lag. For lag $p = 1$, we can use $\{\hat{f}_t\}_{t=1}^T$ to represent time-variant shape functions
 for convenience. We can make dynamic forecasts for the next h windows. One simple forecasting
 method is to use $\hat{W}_{T,k}$ to make forecast for $\hat{\mathbf{X}}_{T+h}$, that is

$$\hat{\mathbf{X}}_{T+h} = \sum_{k=1}^p \hat{W}_{T,k} \mathbf{X}_{T+h-k}.$$

The alternative method is to retrain the new forecast to obtain the latest shape functions as well as
 spatial weight matrix. Given long-term forecast window L , we first make short-term forecast for h
 steps

$$\hat{\mathbf{X}}_{T+h} = \sum_{k=1}^p \hat{W}_{T+h,k} \mathbf{X}_{T+h-k},$$

147 where $h = 1, 2, \dots$ and $\hat{W}_{T+h,k}$ is estimated by training forecast value of $\hat{\mathbf{X}}_{T+h-k}$. We repeat the
 148 process until L steps in total have been predicted.

149 We summarize the whole process of our model when making spatio-temporal forecasts.

- 150 Step 1 Given the observation $\{\mathbf{X}_t\}_{t=1}^T$ and its coordinates, calculate all distance pairs among all
 151 locations, denoted by $\{d_{i_1, j_1}, \dots, d_{i_N, j_N}\}$.
- 152 Step 2 Calculate $\{\frac{1}{m}, \frac{2}{m}, \dots, 1\}$ quantile levels and obtain corresponding distance quantile value
 153 $\{d_{(1)}, d_{(2)}, \dots, d_{(m)}\}$.
- 154 Step 3 Determine the shape constraints and construct corresponding basis functions. Specify the
 155 time lag p .
- 156 Step 4 Train the model according to the illustration of Figure 1.

157 4 Experiment

158 In order to assess our model in stationary and non-stationary scenarios, we synthesize data. Then,
 159 we apply our model to make some comparisons. On the one hand, we need to evaluate how
 160 the estimated shape functions look and assess their similarity and accuracy. On the other hand,
 161 our model can make spatio-temporal forecasting after estimating for spatial weight matrix. The
 162 basic idea for completing the two goals is to set up the expected shape function and compare
 163 estimated parameters with the real one. Next, we assess the forecasting performance with baseline
 164 models. Codes and data for replicating our experiments are anonymously published at <https://anonymous.4open.science/r/STVAR-F16E/>.
 165

166 4.1 Simulation for stationary model

Here, we synthesize 100 stationary spatio-temporal data sets. The spatial domain consists of
 30 locations and their coordinates can be found at <https://anonymous.4open.science/r/STVAR-F16E/>. For each location, we observe 500 values. The observation is generated from
 the stationary model $X_t = \sum_{k=1}^p W_k X_{t-k} + \epsilon_t$, where ϵ_t is randomly generated from the standard
 normal distribution. The next step is to construct random spatial weight matrices for each synthesized
 data set. The shape functions are set to be decreasing, and we set them as a logarithmic function:

$$\alpha(-\log(d+1) + \log(170)),$$

167 where α is randomly generated from uniform distribution [0.05,0.06] but kept to be fixed for each
 168 simulated data set. We use $d+1$ to avoid zero value. This setting can make the real shape function

169 decrease and make it equal to zero when $d = 169$. The stationary model can iteratively generate the
 170 \mathbf{X}_t given initial value \mathbf{X}_0 , where \mathbf{X}_0 is randomly generated from a uniform distribution with bounds
 171 $[-0.01, 0.01]$. The time lags are set as $p = 1$.

172 **Estimation for shape functions.** In Figure 3,
 173 the estimated shape function is presented in red,
 174 while the real shape function is presented in blue.
 175 It can be seen that the estimated shape function
 176 can capture the trend of the real shape function.

177 **Training details.** The first 300 steps are used as
 178 training data, saving the last 200 steps for eval-
 179 uation. We train all models for 100 epochs with
 180 Adam optimizer (5) and a learning rate of 0.01.
 181 The process involves parallel training across 10
 182 CPUs. We select 100 quantile levels, and thus
 183 100 basis functions $g_i(d)$ were generated as the
 184 inputs for the model.

185 **Assessment for forecasting.** We assess the fore-
 186 casting performance for the stationary model
 187 with baseline models. As introduced in the liter-
 188 ature review, the baseline models are selected from the VAR model(21), the spatial panel data(SPE)
 189 model that applied pre-specified spatial weigh matrix (28), graph-based models (20; 37) and
 190 convolution-based models (15; 23). The error metrics are mean absolute error and root mean squared
 191 error defined by $\frac{1}{Nn} \sum_{j=1}^N \sum_{i=1}^n \frac{\sum_{t=T}^{T+h} |\hat{X}_{it}^{(j)} - X_{it}^{(j)}|}{h}$, $\frac{1}{Nn} \sum_{j=1}^N \sum_{i=1}^n \sqrt{\frac{1}{h} \sum_{t=T}^{T+h} (X_{it}^{(j)} - \hat{X}_{it}^{(j)})^2}$,
 192 respectively. The Table 1 shows the six baseline models with the proposed model. As totally we
 193 have 100 synthesised data sets, $X_{it}^{(j)}$ and $\hat{X}_{it}^{(j)}$ denote $i - th$ variable in $j - th$ data sets. $n = 30$ is
 194 the number of locations and $N = 100$ is the number of synthesised data. We conducted one-step
 195 forecasting for the next 200 observations.

196 Compared with baseline models, the proposed model performs better under the metric MAE and
 197 RMSE. The proposed method outperforms the closest competing method, DC-RNN, by 10%.

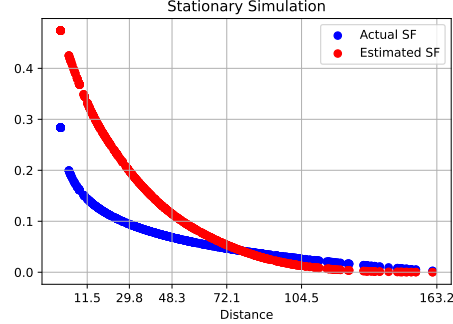


Figure 3: The sample of estimated shape function. Distances are shown every 20^{th} quantile.

198 4.2 Experiments for non-stationary model

We conduct a simulation for the non-stationary model with time lag $p = 1$ and synthesize 100 data sets using a similar approach to the stationary model simulation. The initial value \mathbf{X}_0 and ϵ_t are generated from a uniform and normal distribution respectively. The locations of observations are the same as those in the stationary model simulation. In order to construct W_t , the time-varying shape functions are created under the decreased constraint. The shape function at time t is constructed as

$$\alpha_t(-\log(d+1) + \log(170)),$$

199 where α_t controls the level of value at each time t . ϵ_t is generated from a normal distribution. \mathbf{X}_0 is
 200 generated from a uniform distribution with bound $[-0.001, 0.001]$.

201 **Shape functions settings and estimation.** The shape functions are set as time-variant, as they can
 202 simulate the non-stationary process across time. We specified α_0 at $t = 0$ from uniform distribution
 203 $[1 \times 10^{-4}, 2 \times 10^{-4}]$ and then make an interpolation from α_0 to α_{500} . The total length for every
 204 location is 500 and we set $\alpha_{500} = 10 \times \alpha_0$. For example, generally if $\alpha_0 = 0.0001$, we have
 205 $\alpha_t = 0.0001(1 - \frac{t}{T}) + 0.001 \frac{t}{T}$, where $T = 500$. This setting guarantee that shape functions vary
 206 from lower level to higher level. The larger α_t is, the more larger distance-based effects they have.
 207 Thus, the corresponding spatial weight matrix consists of dynamic shape functions and can reflect the
 208 non-stationary dependence among each site. We present the estimated shape functions in Figure 4
 209 and compare them with the real ones.

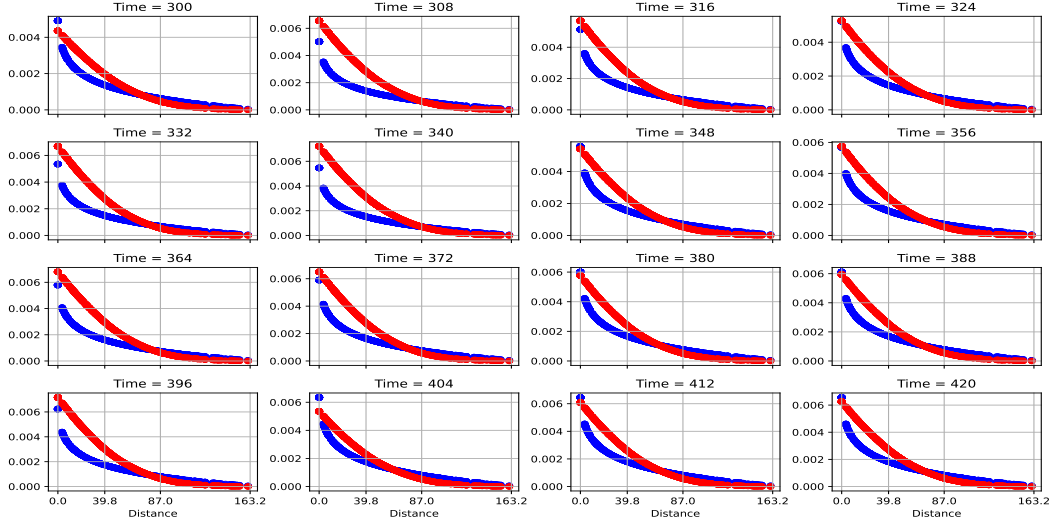


Figure 4: The sample of estimated shape function for the 120 testing time steps. Distances are shown every 40th quantile.

210 **Training details.** Similar to the stationary simulation, the train-test split is 300 – 200 over the data
 211 size of 500. However, we train all models for 100 epochs with Adam optimizer (5) at a learning rate
 212 of 0.001. We train models in parallel across 10 CPUs.

213 **Forecasting performance.** The forecasting performance is assessed by the same metrics used in the
 214 previous simulation for the stationary case. We made a one-step forecast by our model. As for the
 215 baseline models, we adjusted their published code accordingly. The results show that the proposed
 216 model can still capture non-stationary processes compared with baseline models. The proposed
 217 method outperforms the other competing methods. The error metric is shown in Table 1.

Table 1: The error metrics with baseline models for simulation.

Methods	Stationary Simulation		Non-stationary Simulation	
	MAE	RMSE	MAE	RMSE
VAR	2.9611 ± 1.8573	3.2588 ± 1.8077	2.4426 ± 1.2285	2.7676 ± 1.2015
SPM	1.8850 ± 0.6348	1.8671 ± 0.6778	2.1918 ± 0.7350	2.2161 ± 0.6876
DC-RNN	0.8960 ± 0.0370	1.1168 ± 0.0426	0.9017 ± 0.0358	1.1328 ± 0.0463
FC-GAGA	2.5425 ± 0.2965	3.1066 ± 0.3633	1.0270 ± 0.0080	1.2939 ± 0.0120
GMAN	1.6806 ± 0.1491	1.9293 ± 0.1483	1.5714 ± 0.1104	1.8608 ± 0.1155
ConvLSTM	2.9495 ± 0.2980	3.2509 ± 0.2887	2.2478 ± 0.2295	2.5469 ± 0.2324
ESTF	0.7997 ± 0.0015	1.0017 ± 0.0016	0.8075 ± 0.0016	1.0112 ± 0.0020

218 4.3 Real case studies

219 **Air quality data.** We apply our model to air quality data, which records air quality in California over
 220 2021¹. The daily mean of PM 2.5 is recorded across 172 sites.

221 We obtain the first 200 steps for training and perform forecasting for the next 165 steps. All models
 222 are trained for 100 epochs using Adam optimizer (5), at a learning rate of 0.01 and batch size of
 223 50. We present the estimated time-variant shape functions in supplemental file. The value of shape
 224 functions decays to zero at around 5.926, which is 80% quantile in the sample of distance pairs. In
 225 other words, the distance-based effects decay to zero at a distance equal or larger than 5.926. Our

¹<https://www.epa.gov/outdoor-air-quality-data/download-daily-data>

226 model has ideal performance with low time consumption compared with baseline models. We put
 227 detailed forecasting results of simulation and real cases in a supplemental file.

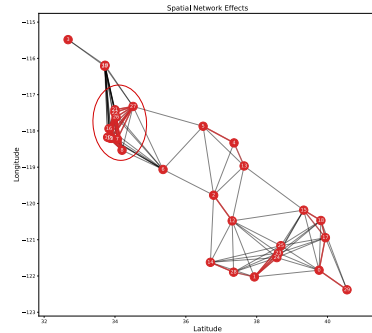
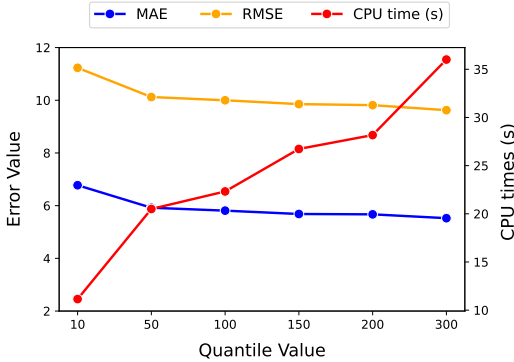


Figure 5: Comparing efficiency vs. performance trade-off at different quantile values. Figure 6: The significant distance-based effect among all 30 locations.

228 The result is shown in Table 2. The ESTF performs best in terms of RMSE, while the DC-RNN
 229 method performs best in terms of MAE. For the computational time, the ESTF method is significantly
 230 faster than most machine learning methods, and only takes around 1/10 time of DC-RNN. In Figure 5,
 231 MAE, RMSE, and time are presented with different numbers of m . As m increases, the computational
 232 time increases while both MAE and RMSE decrease. There is a significant increase in the forecasting
 233 performance when m increases from 10 to 50. For $m > 50$, the forecasting performance does not
 234 increase much as m increases.

235 One key advantage of the ESTF method is that we can make an explicit distance-based explanation
 236 for our dataset. Figure 6 shows the distance-based effects at time $t = 9$. We only present the effects
 237 using a threshold to obtain a more concise visualization. The estimated shape function \hat{f}_9 ranges
 238 from 0 to 9.8 and we set 5 as the threshold. The red line indicates the value of the shape function
 239 larger than 7, while the gray line indicates the value between 5 and 7. Figure 6 shows how any two
 240 locations interact and measure the distance-based effects quantitatively. For example, air quality
 241 monitoring sites around the Greater Los Angeles (red circle in Figure 6) area have a strong spatial
 242 interaction with each other, such as node 7 and node 8.

243 4.4 SO_2 data

244 Texas is the second largest manufacturing state in the USA and prediction for SO_2 is critical task for
 245 researchers. The data² records daily SO_2 at 31 locations in 2021. More detailed spatial information
 246 can be found in the supplemental file. The numeric result is listed in Table 2.

Table 2: The error metrics with baseline models for real case study. Clock time (in seconds) for real case study is recorded when training each model for 100 epochs on a single CPU.

Methods	Air quality data				SO_2 data			
	MAE	RMSE	Training time (s)	Inference Time (s)	MAE	RMSE	Training Time (s)	Inference Time (s)
VAR	16.9844	22.3410	3.56	0.04	6.2705	9.1388	3.330	0.016
SPM	8.4547	13.8262	0.31	0.03	7.1453	9.1086	0.143	0.027
DC-RNN	4.7157	9.3873	203	1.211	3.5094	6.8681	264.215	1.366
FC-GAGA	7.8671	18.1870	181	2.759	4.5976	7.7528	169.425	2.889
GMAN	12.5268	17.3817	140	1.823	4.1099	7.4806	172.016	1.581
ConvLSTM	12.6292	17.9149	53	1.940	4.1445	8.0688	96.233	1.656
ESTF	5.2237	9.2169	22	1.625	4.2966	6.8307	31.050	1.868

247 Similar conclusions can be drawn in SO_2 data as that of air quality data. The ESTF model performs
 248 best under the RMSE metric, while DC-RNN is best in the MAE metric. In terms of training time,

²<https://www.epa.gov/outdoor-air-quality-data/download-daily-data>

249 the proposed ESTF method costs around 10% of that of DC-RNN. For the inference time, ESTF and
 250 DC-RNN are comparable.

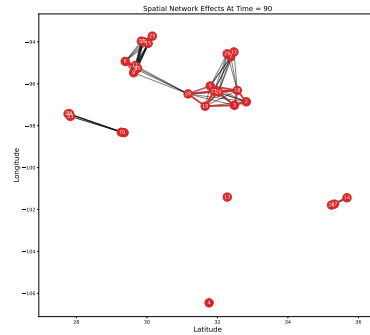
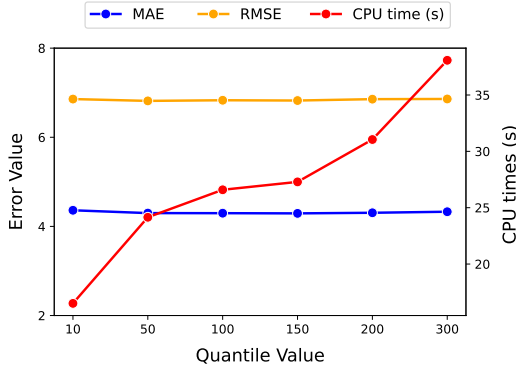


Figure 7: Comparing efficiency vs. performance among all 31 locations at $t = 90$ trade-off at different quantile values. Figure 8: The significant distance-based effect

251 The efficiency analysis and performance at different quantiles are shown in Figure 7. Together with
 252 Figure 5, we can see that the increasing number of basis functions does not have much improvement
 253 when the number of basis functions is larger than 50, while the training time increases as the number
 254 of basis functions increases. The spatial distribution at time $t = 90$ is presented in Figure 8 where
 255 coordinates are denoted by latitude and longitude. Two significant clusters, representing Houston
 256 and Dallas respectively, have the strongest distance-based effect. It quantitatively shows how these
 257 neighbors affect each other. Counties around Dallas-Fort Worth metropolitan area show strong
 258 interaction, which should be noted by environmental policy-makers. More detailed results are
 259 presented in the supplemental file.

260 5 Discussion

261 This paper applies learnable shape functions to capture distance-based effects. It can model dynamic
 262 spatial dependence for stationary and non-stationary spatio-temporal data based on their distance.
 263 The model does not have the limitations of classical statistical spatial models and provides a more
 264 explanatory model than usual deep learning methods. Furthermore, some spatio-temporal data, such
 265 as temperature for sea surface and air quality monitoring data, usually viewed as collected from the
 266 continuous field, are more suitable for the proposed models since these kinds of data follow the basic
 267 rule that variability between two locations is significantly affected by their distance. However, some
 268 spatio-temporal data, such as traffic flow or some biology data, do not follow the rule. As a result, the
 269 spatial dependence may rely on road structure or biological mechanisms instead of distance. It is
 270 worth researching such data by considering graph structure when estimating spatial weight matrix. In
 271 addition, we can develop spatio-temporal causal inference based on the ESTF model. Grander causal
 272 analysis can be done by fitting the first-order VAR model (24). The estimation of the coefficients
 273 matrix of the VAR model attracts researchers' interest as it can be treated as a causal transition
 274 matrix. In the causal inference community, lots of work have been conducted on the VAR model
 275 (8; 9). However, there is a lack of research on causal inference under the spatio-temporal process.
 276 The quantitative distance-based effects in ESTF can be further researched and extended to develop a
 277 spatio-temporal causal model.

278 References

- 279 [1] ANSELIN, L. *Spatial econometrics: methods and models*, vol. 4. Springer Science & Business
 280 Media, 1988.
- 281 [2] CASTRUCCIO, S., AND GENTON, M. G. Principles for statistical inference on big spatio-
 282 temporal data from climate models. *Statistics & Probability Letters* 136 (2018), 92–96.

- 283 [3] CHEN, Y., AND SAMWORTH, R. J. Generalized additive and index models with shape
284 constraints. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 78, 4
285 (2016), 729–754.
- 286 [4] CLIFF, A. Spatial autocorrelation: Technical report.
- 287 [5] DIEDERIK, K., JIMMY, B., ET AL. Adam: A method for stochastic optimization. *arXiv preprint*
288 *arXiv:1412.6980* (2014), 273–297.
- 289 [6] DOU, B., PARRELLA, M. L., AND YAO, Q. Generalized Yule–Walker estimation for spatio-
290 temporal models with unknown diagonal coefficients. *Journal of Econometrics* 194, 2 (2016),
291 369–382.
- 292 [7] ELHORST, J. P. Spatial panel data models. In *Spatial econometrics*. Springer, 2014, pp. 37–93.
- 293 [8] GEIGER, P., ZHANG, K., SCHOELKOPF, B., GONG, M., AND JANZING, D. Causal inference
294 by identification of vector autoregressive processes with hidden components. In *International*
295 *Conference on Machine Learning* (2015), PMLR, pp. 1917–1925.
- 296 [9] GONG, M., ZHANG, K., SCHÖLKOPF, B., GLYMOUR, C., AND TAO, D. Causal discovery
297 from temporally aggregated time series. In *Uncertainty in artificial intelligence: proceedings of*
298 *the... conference. Conference on Uncertainty in Artificial Intelligence* (2017), vol. 2017, NIH
299 Public Access.
- 300 [10] JAIN, A., ZAMIR, A. R., SAVARESE, S., AND SAXENA, A. Structural-rnn: Deep learning on
301 spatio-temporal graphs. In *Proceedings of the ieee conference on computer vision and pattern*
302 *recognition* (2016), pp. 5308–5317.
- 303 [11] KELEJIAN, H. H., AND PRUCHA, I. R. A generalized spatial two-stage least squares procedure
304 for estimating a spatial autoregressive model with autoregressive disturbances. *The Journal of*
305 *Real Estate Finance and Economics* 17, 1 (1998), 99–121.
- 306 [12] LEE, L.-F. Asymptotic distributions of quasi-maximum likelihood estimators for spatial
307 autoregressive models. *Econometrica* 72, 6 (2004), 1899–1925.
- 308 [13] LEE, L.-F., AND YU, J. Some recent developments in spatial panel data models. *Regional*
309 *Science and Urban Economics* 40, 5 (2010), 255–271.
- 310 [14] LI, M., AND ZHU, Z. Spatial-temporal fusion graph neural networks for traffic flow forecasting.
311 In *Proceedings of the AAAI conference on artificial intelligence* (2021), vol. 35, pp. 4189–4196.
- 312 [15] LI, Y., YU, R., SHAHABI, C., AND LIU, Y. Diffusion convolutional recurrent neural network:
313 Data-driven traffic forecasting. *arXiv preprint arXiv:1707.01926* (2017).
- 314 [16] LIN, Y., MAGO, N., GAO, Y., LI, Y., CHIANG, Y.-Y., SHAHABI, C., AND AMBITE, J. L.
315 Exploiting spatiotemporal patterns for accurate air quality forecasting using deep learning. In
316 *Proceedings of the 26th ACM SIGSPATIAL international conference on advances in geographic*
317 *information systems* (2018), pp. 359–368.
- 318 [17] MA, Y., GUO, S., AND WANG, H. Sparse spatio-temporal autoregressions by profiling and
319 bagging. *Journal of Econometrics* (2021).
- 320 [18] MILLER, H. J. Tobler’s first law and spatial analysis. *Annals of the association of American*
321 *geographers* 94, 2 (2004), 284–289.
- 322 [19] MOKBEL, M. F., XIONG, X., HAMMAD, M. A., AND AREF, W. G. Continuous query
323 processing of spatio-temporal data streams in place. *GeoInformatica* 9, 4 (2005), 343–365.
- 324 [20] ORESHKIN, B. N., AMINI, A., COYLE, L., AND COATES, M. J. FC-GAGA: Fully connected
325 gated graph architecture for spatio-temporal traffic forecasting. In *Proc. AAAI Conf. Artificial*
326 *Intell* (2021).
- 327 [21] QIAN, G., TORDESILLAS, A., AND ZHENG, H. Landslide forecast by time series modeling
328 and analysis of high-dimensional and non-stationary ground motion data. *Forecasting* 3, 4
329 (2021), 850–867.

- 330 [22] QU, X., LEE, L.-F., AND YU, J. QML estimation of spatial dynamic panel data models with
331 endogenous time varying spatial weights matrices. *Journal of Econometrics* 197, 2 (2017),
332 173–201.
- 333 [23] SHI, X., CHEN, Z., WANG, H., YEUNG, D.-Y., WONG, W.-K., AND WOO, W.-C. Convolutional
334 lstm network: A machine learning approach for precipitation nowcasting. *Advances in
335 neural information processing systems* 28 (2015).
- 336 [24] SIMS, C. A. Macroeconomics and reality. *Econometrica: journal of the Econometric Society*
337 (1980), 1–48.
- 338 [25] SU, L. Semiparametric GMM estimation of spatial autoregressive models. *Journal of Econo-
339 metrics* 167, 2 (2012), 543–560.
- 340 [26] TOBLER, W. R. A computer movie simulating urban growth in the detroit region. *Economic
341 geography* 46, sup1 (1970), 234–240.
- 342 [27] WANG, D., AND CHENG, T. A spatio-temporal data model for activity-based transport demand
343 modelling. *International Journal of Geographical Information Science* 15, 6 (2001), 561–585.
- 344 [28] WANG, H., QIAN, G., AND TORDESILLAS, A. Modeling big spatio-temporal geo-hazards data
345 for forecasting by error-correction cointegration and dimension-reduction. *Spatial Statistics* 36
346 (2020), 100432.
- 347 [29] WANG, S., CAO, J., AND YU, P. Deep learning for spatio-temporal data mining: A survey.
348 *IEEE transactions on knowledge and data engineering* (2020).
- 349 [30] WANG, X., CHEN, C., MIN, Y., HE, J., YANG, B., AND ZHANG, Y. Efficient metropolitan
350 traffic prediction based on graph recurrent neural network. *arXiv preprint arXiv:1811.00740*
351 (2018).
- 352 [31] WANG, Y., LONG, M., WANG, J., GAO, Z., AND YU, P. S. PredRNN: Recurrent neural
353 networks for predictive learning using spatio-temporal LSTMs. In *Advances in Neural Informa-
354 tion Processing Systems* (2017), I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus,
355 S. Vishwanathan, and R. Garnett, Eds., vol. 30, Curran Associates, Inc.
- 356 [32] WU, Y., AND TAN, H. Short-term traffic flow forecasting with spatial-temporal correlation in a
357 hybrid deep learning framework. *arXiv preprint arXiv:1612.01022* (2016).
- 358 [33] YANG, S., MA, W., PI, X., AND QIAN, S. A deep learning approach to real-time parking
359 occupancy prediction in transportation networks incorporating multiple spatio-temporal data
360 sources. *Transportation Research Part C: Emerging Technologies* 107 (2019), 248–265.
- 361 [34] YU, B., YIN, H., AND ZHU, Z. Spatio-temporal graph convolutional networks: A deep
362 learning framework for traffic forecasting. *arXiv preprint arXiv:1709.04875* (2017).
- 363 [35] YU, J., DE JONG, R., AND LEE, L.-F. Quasi-maximum likelihood estimators for spatial
364 dynamic panel data with fixed effects when both n and t are large. *Journal of Econometrics* 146,
365 1 (2008), 118–134.
- 366 [36] YUAN, Z., ZHOU, X., AND YANG, T. Hetero-convlstm: A deep learning approach to traffic
367 accident prediction on heterogeneous spatio-temporal data. In *Proceedings of the 24th ACM
368 SIGKDD International Conference on Knowledge Discovery & Data Mining* (2018), pp. 984–
369 992.
- 370 [37] ZHENG, C., FAN, X., WANG, C., AND QI, J. Gman: A graph multi-attention network for
371 traffic prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2020),
372 vol. 34, pp. 1234–1241.
- 373 [38] ZHOU, S., BONDELL, H., TORDESILLAS, A., RUBINSTEIN, B. I., AND BAILEY, J. Early
374 identification of an impending rockslide location via a spatially-aided gaussian mixture model.
375 *The Annals of Applied Statistics* 14, 2 (2020), 977–992.

376 **Checklist**

- 377 1. For all authors...
- 378 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
379 contributions and scope? [Yes]
- 380 (b) Did you describe the limitations of your work? [Yes] We describe limitations in
381 discussion section
- 382 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 383 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
384 them? [Yes]
- 385 2. If you are including theoretical results...
- 386 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 387 (b) Did you include complete proofs of all theoretical results? [N/A]
- 388 3. If you ran experiments...
- 389 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
390 mental results (either in the supplemental material or as a URL)? [Yes] We upload data
391 and code to github
- 392 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
393 were chosen)? [Yes] Please check training details in simulation and real case section
- 394 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
395 ments multiple times)? [Yes]
- 396 (d) Did you include the total amount of compute and the type of resources used (e.g.,
397 type of GPUs, internal cluster, or cloud provider)? [Yes] They are included in training
398 details
- 399 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 400 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 401 (b) Did you mention the license of the assets? [Yes]
- 402 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 403 (d) Did you discuss whether and how consent was obtained from people whose data you're
404 using/curating? [N/A]
- 405 (e) Did you discuss whether the data you are using/curating contains personally identifiable
406 information or offensive content? [N/A]
- 407 5. If you used crowdsourcing or conducted research with human subjects...
- 408 (a) Did you include the full text of instructions given to participants and screenshots, if
409 applicable? [N/A]
- 410 (b) Did you describe any potential participant risks, with links to Institutional Review
411 Board (IRB) approvals, if applicable? [N/A]
- 412 (c) Did you include the estimated hourly wage paid to participants and the total amount
413 spent on participant compensation? [N/A]