Spatial Functional Linear Model and Its Estimation Method

Tingting Huang¹, Gilbert Saporta², Huiwen Wang^{1,3}, Shanshan Wang^{1,4}

 ¹School of Economics and Management,Beihang University,Beijing,China
 ²CEDRIC, CNAM, Paris
 ³Beijing Advanced Innovation Center for Big Data and Brain Computing, Beihang University, Beijing ⁴Beijing Key Laboratory of Emergence Support Simulation Technologies for City Operations,Beijing,China

Dec., 2019

Spatial Functional Linear Model

T. Huang, G. Saporta, H. Wang and S. Wang

- Functional linear model and motivation
- 2 Spatial functional linear model
- Stimation method
- Simulation study
- Weather data
- Onclusion

What is functional data?

Functional data is a type of data recorded at a high frequency , and almost can be regarded as continuously observed. trajectory data, weather data and stock data. functional data are inherently

- infinite-dimensional
- rich in information, derivative
- biology, medical sciences, econometrics

Functional linear model (FLM) becomes more and more popular owing to its adaptability to problems which are hard to deal with in the framework of scalar and vector observations.

Cai and Hall (2006); Hall and Horowitz (2007); Ferraty et al. (2013)

$$Y = \alpha + \int_{\Gamma} X(t)\beta(t)dt + \epsilon, \qquad (1)$$

- Y : a scalar response variable
- X(t): a second-order stochastic process on a compact interval Γ , E(X(t)) = 0, $E(\int_{\Gamma} X^2(t) dt) < \infty$
- $\beta(t)$: unknown slope function
- *ϵ* : random error term independent of X(t) with zero mean and finite variance

Spatial Functional Linear Model

Extensions of FLM have been studied to address specific problems.

- James (2002)put forward functional logistic regression and functional censored regression
- Aneiros-Prez and Vieu (2006) constructed a semi-functional partial linear model
- Ferraty et al. (2013) generalized the FLM to functional projection pursuit regression that allows for more interpretability
- Liu et al. (2017) presented a functional linear mixed model

Secondly, relation information among individuals with advance of information technology can be easily collected.

And as one typical relational information, the network-structured data are becoming increasing available.

Secondly, relation information among individuals with advance of information technology can be easily collected.

And as one typical relational information, the network-structured data are becoming increasing available.

• If we use FLM or the aforementioned variations of FLM to model this kind of data

Secondly, relation information among individuals with advance of information technology can be easily collected.

And as one typical relational information, the network-structured data are becoming increasing available.

- If we use FLM or the aforementioned variations of FLM to model this kind of data
 - the information contained in the data may not be fully exploited

Secondly, relation information among individuals with advance of information technology can be easily collected.

And as one typical relational information, the network-structured data are becoming increasing available.

- If we use FLM or the aforementioned variations of FLM to model this kind of data
 - the information contained in the data may not be fully exploited
 - The inference may be misleading when ignoring such network effects

Weather data of temperature and precipitation in cities of China

- Data source : China Meteorological Yearbook
- Description : monthly temperature and precipitation in 34 major cities of China between 2005 and 2007
- Problem : effect of temperature on the precipitation over three years

Y is the logarithm of mean annual precipitation X(t) is the mean monthly temperature

1.4 motivating example





key cities

monthly temperature in a year

We adopt Moran I statistic to test whether there exists spatial autocorrelation between the responses or not

the original precipitation Moran I statistic = 0.7P value < 0.000 001 the residuals of FLM Moran I statistic = 0.5P value < 0.000 001

we can observe that the responses as well as the residuals of FLM showed a significance spatial autocorrelation, which means that FLM may be no longer appropriate for such data.

1.4 motivating example



moran I plot of precipitation

Moran I plot of residuals of FLM

When the predictor is scalar, the spatial autoregressive model (SAR) is the popularly used model to accommodate the dependence brought by the network structure.

$$\mathbf{y} = \alpha \boldsymbol{\tau}_n + \rho \mathbf{W} \mathbf{y} + \mathbf{x} \boldsymbol{\beta} + \boldsymbol{\epsilon}$$
(2)

- a spatial weight matrix is adopted to denote adjacent relations
- an unknown spatial autoregressive parameter is used to reflect the strength of neighboring effects

Borrowing thoughts from SAR model, we propose a new model which is Spatial Functional Linear Model (SFLM).

Model Assumption

- Following Qu and Lee (2015), we consider the spatial process located on a unevenly spaced lattice D ⊆ R^d, d ≥ 1.
- We observe $\{(x_i(t), y_i)\}_{i=1}^n$ from *n* spatial units on *D*.
- x_i(t)'s are independent and identically distributed (i.i.d.) samples from X(t), where X(t) is a square integrable second-order stochastic process defined on a compact set Γ with E(X(t)) = 0, and E(∫_Γ X²(t)dt) < ∞.

$$\mathbf{y} = \alpha \boldsymbol{\tau}_n + \rho \mathbf{W} \mathbf{y} + \int_0^1 \mathbf{x}(t) \beta(t) dt + \boldsymbol{\epsilon}$$
(3)

2. spatial functional linear model

$$\mathbf{y} = lpha \boldsymbol{\tau}_n +
ho \mathbf{W} \mathbf{y} + \int_0^1 \mathbf{x}(t) \beta(t) dt + \epsilon$$

- $\alpha \boldsymbol{\tau}_n$ is the intercept, α denotes a scalar parameter
- ρ is the unknown spatial autocorrelation parameter takes value from [0, 1)
- $\mathbf{W} = (w_{ii'})_{n \times n}$ is a pre-specified spatial weight matrix
- β(t) is a square integrable coefficient function defined on [0, 1]
- ϵ is the noise term independent of $\mathbf{x}(\mathbf{t})$, follows a multivariate normal distribution with zero mean and a constant diagonal covariance matrix $\sigma^2 \mathbf{I}_n$

t

notes of spatial functional linear model

- the spatial weight matrix W is exogenous
- ρ is a scale parameter measuring the power of the neighbor impact
- the SFLM can degenerate into the following two models
 - functional linear model (FLM) when $\rho = 0$
 - $\bullet\,$ spatial autoregressive model (SAR) when x(t) is free of

we can reformulate model (3) as the following equivalent expression as well,

$$\mathbf{y} = (\mathbf{I}_n - \rho \mathbf{W})^{-1} (\alpha \boldsymbol{\tau}_n + \int_0^1 \mathbf{x}(\mathbf{t}) \beta(t) dt + \epsilon) \qquad (4)$$

y_i is influenced by its neighbors' covariates x_{i'}(t)s, i' ≠ i
Gauss-Markov assumption is violated, the least square estimator (OLS) is not adequate

- Basis Expansion based on FPCA
- In MLE for truncated SFLM

Basis Expansion based on FPCA

$$\begin{split} & \mathcal{K}(s,t) = \mathsf{Cov}(X(t),X(s)). \\ & \mathcal{K}(s,t) = \sum_{j=1}^{\infty} k_j \varphi_j(s) \varphi_j(t) \ , \quad k_1 > k_2 > \cdots > 0. \\ & X(t) = \sum_{j=1}^{\infty} a_j \varphi_j(t). \end{split}$$

For the observation $\{y_i, x_i(t)\}_{i=1}^n$,

$$\hat{K}(s,t) = rac{1}{n} \sum_{i=1}^{n} x_i(s) x_i(t) - \bar{x}(s) \bar{x}(t) , \ \bar{x}(t) = rac{1}{n} \sum_{i=1}^{n} x_i(t)$$

Spatial Functional Linear Model

T. Huang, G. Saporta, H. Wang and S. Wang

3. estimation method

Basis Expansion based on FPCA

Also we can obtain that

$$\hat{K}(s,t) = \sum_{j=1}^{\infty} \hat{k}_j \hat{arphi}_j(s) \hat{arphi}_j(t)$$

$$\hat{a}_{ij}=\int_0^1 x_i(t)\hat{arphi}_j(t)\;,\;x_i(t)=\sum_{j=1}^\infty \hat{a}_{ij}\hat{arphi}_j(t)$$

$$b_j = \int_0^1 eta(t) \hat{arphi}_j(t) dt \;,\; eta(t) = \sum_{j=1}^\infty b_j \hat{arphi}_j(t)$$

$$\mathbf{y} = \alpha \boldsymbol{\tau}_n + \rho \mathbf{W} \mathbf{y} + \sum_{j=1}^{\infty} \hat{\mathbf{a}}_j b_j + \boldsymbol{\epsilon}, \qquad (5)$$

Spatial Functional Linear Model

Basis Expansion based on FPCA

Using the criterion, percentage of variance explained (PVE) for predictors to choose number of PCs. The truncated model is

$$\mathbf{y} \approx \alpha \boldsymbol{\tau}_n + \rho \mathbf{W} \mathbf{y} + \sum_{j=1}^m \hat{\mathbf{a}}_j b_j + \boldsymbol{\epsilon}.$$
 (6)

Based on estimated \hat{b}_j , we can reconstruct $\hat{eta}(t)$ by

$$\hat{\beta}(t) = \sum_{j=1}^{m} \hat{b}_j \hat{\varphi}_j(t)$$
(7)

MLE for truncated SFLM

Defining $\mathbf{A} = (\hat{a}_{ij})_{n \times m}$, $\mathbf{b} = (b_1, b_2, \cdots, b_m)'$, $\mathbf{Z} = (\tau_n, \mathbf{A})$ and $\boldsymbol{\delta} = (\alpha, \mathbf{b}')'$, we can write the truncated model (6) as

$$\mathbf{y} \approx \rho \mathbf{W} \mathbf{y} + \mathbf{Z} \boldsymbol{\delta} + \boldsymbol{\epsilon}.$$
 (8)

Since model (8) is similar to the popular SAR model, we adopt the maximum likelihood estimation method to get parameters.

- Represent the functional predictor and slope function by functional principle component basis, then the estimation procedure is simplified as the remaining process is similar to the estimation problem of a SAR model.
- Get estimators of unknown parameters in the truncated SFLM obtained from Step 1.
- Get the estimator of β(t) in SFLM. The slope function is constructed with the FPC basis in Step 1 and the estimated coefficients in Step 2. Other estimators are directly obtained from Step 2.

- As for spatial scenario, we adopt the rook matrix by randomly apportioning *n* agents on a regular square grid of cells.
- And for the functional part in Model (2), we take the same form as functions in FLM by Hall and Horowitz (2007).

$$Y = \rho WY + \int_0^1 X(t)\beta(t)dt + 0.5\epsilon$$
$$X(t) = \sum_{j=1}^{50} a_j Z_j \varphi_j(t) \qquad \beta(t) = \sum_{j=1}^{50} b_j \varphi_j(t)$$
$$a_j = (-1)^{j+1} j^{-\gamma/2} \quad Z_j \sim U[-\sqrt{3}, \sqrt{3}] \qquad \epsilon \sim N[0, 1]$$
$$b_1 = 0.3 \qquad b_j = 4(-1)^{j+1} j^{-2}, j \ge 2 \qquad \varphi_j(t) = \sqrt{2}cos(j\pi t)$$
$$\rho = \{0, 0.5, 0.8\} \quad \gamma = \{1.1, 2\} \quad n = 300, 500, 900$$

Spatial Functional Linear Model

			$\gamma = 1.1$			$\gamma = 2$	
ρ	n	bias(sd)	$MSE_1(sd)$	$MSE_2(sd)$	bias(sd)	$MSE_1(sd)$	$MSE_2(sd)$
0	300	-0.0051 $_{(0.050)}$	$\underset{(0.007)}{0.0203}$	$\underset{(0.007)}{0.0203}$	-0.0086 $_{(0.063)}$	$\underset{(0.024)}{0.1171}$	$\underset{(0.024)}{0.1171}$
	500	-0.0003 $_{(0.043)}$	$\substack{0.0087\(0.003)}$	$\substack{0.0087\(0.003)}$	-0.0020 $_{(0.046)}$	$\underset{(0.011)}{0.0691}$	$\underset{(0.011)}{0.0691}$
	900	-0.0024 $_{(0.029)}$	$\underset{(0.001)}{0.0034}$	$\underset{(0.001)}{0.0034}$	-0.0024 $_{(0.029)}$	$\underset{(0.001)}{0.0034}$	$\underset{(0.001)}{0.0034}$
0.5	300	-0.0062 $_{(0.046)}$	$\underset{(0.007)}{0.0201}$	$\substack{0.0267\ (0.010)}$	-0.0068 $_{(0.052)}$	$\underset{(0.023)}{0.1170}$	$\underset{(0.023)}{0.1198}$
	500	-0.0016 $_{(0.034)}$	$\underset{(0.003)}{0.0085}$	$\underset{(0.004)}{0.0114}$	-0.0037 $_{(0.040)}$	$\underset{(0.010)}{0.0689}$	$\underset{\scriptscriptstyle(0.010)}{0.0702}$
	900	-0.0034 $_{(0.024)}$	$\underset{(0.001)}{0.0034}$	$\substack{0.0047\(0.001)}$	-0.0046 $_{(0.029)}$	$\underset{(0.005)}{0.0380}$	$\underset{(0.005)}{0.0386}$
0.8	300	-0.0062 $_{(0.026)}$	$\underset{(0.007)}{0.0202}$	$\underset{(0.036)}{0.0836}$	-0.0094	$\underset{(0.023)}{0.1173}$	$\underset{(0.033)}{0.1504}$
	500	-0.0027 $_{(0.020)}$	$\substack{0.0087\(0.003)}$	$\underset{(0.015)}{0.0405}$	-0.0057 $_{(0.025)}$	$\underset{(0.011)}{0.0689}$	$\underset{(0.016)}{0.0876}$
	900	-0.0024 $_{(0.015)}$	$\underset{(0.001)}{0.0033}$	$\underset{(0.006)}{0.0190}$	-0.0040 $_{(0.019)}$	$\underset{(0.004)}{0.0382}$	$\underset{(0.006)}{0.0479}$

$MSE_1 = \frac{1}{n}\sum_{i=1}^n (\widehat{\beta}_{SFLM}(t_i) - \beta(t_i))^2$, $MSE_2 = \frac{1}{n}\sum_{i=1}^n (\widehat{\beta}_{FLM}(t_i) - \beta(t_i))^2$

Spatial Functional Linear Model

T. Huang, G. Saporta, H. Wang and S. Wang



Spatial Functional Linear Model

T. Huang, G. Saporta, H. Wang and S. Wang

- When $\rho = 0$, our proposed method and FPCA based method perform equally well.
- **When** $\rho \neq 0$, our proposed method behaves better than FPCA based method.
- No matter what ρ equals to, MSE₁ of β(t) decreases as sample size increases. And the standard deviation of ρ has a decreasing pattern.
- As the case in Hall and Horowitz (2007), $\beta(t)$ is better estimated with $\gamma = 1.1$ than $\gamma = 2$ when other simulation parameters equal.

We perform the SFLM and FLM on the average of weather data from 2005 to 2007 and *compare their efficiency*. The SFLM is formulated as

$$y_i = \rho \sum_{i \neq i'} w_{ii'} y_{i'} + \int_0^1 x_i(t) \beta(t) dt + \epsilon_i, \qquad (9)$$

where $w_{ii'}$ is the weight between city *i* and *i'*. We also built the FLM as

$$y_i = \int_0^1 x_i(t)\beta(t)dt + \epsilon_i \tag{10}$$

Main steps

- we smoothed the mean monthly temperature over 3 years by Epanechnikov Kernel.
- **2** the spatial weight matrix was formed by nearest 5 neighbors, each neighbor's weight equalling reciprocal of Euclidean distance d(i, i') between cities *i* and *i'*.
- Urumchi's location is far from other cities, we removed its record in weather data.
- We perform the SFLM in model (10) and FLM in model (11) to get estimators.
- Then we apply the fitting result to temperature observations in 2008 to predict annual precipitation.

5. weather data



We conclude that the precipitation is much more influenced by temperature in winter than in other seasons. And under SFLM the precipitation of each city is less affected by temperature during the whole year.

Table 2. The fitting and predicting results

models	$\widehat{ ho}$	fitted error	predicted error	Moran I statistic
FLM	0.66	0.40	0.16	0.50
SFLM	_	0.25	0.13	0.13

We propose a powerful spatial functional linear model which integrates the advantages of FLM in dealing with high dimensional data and SAR model in coping with spatial dependence.

- We propose a powerful spatial functional linear model which integrates the advantages of FLM in dealing with high dimensional data and SAR model in coping with spatial dependence.
- A simple estimation method is developed to obtain estimators

- We propose a powerful spatial functional linear model which integrates the advantages of FLM in dealing with high dimensional data and SAR model in coping with spatial dependence.
- A simple estimation method is developed to obtain estimators
- Our simulation study demonstrate the consistency of the proposed estimators.

- We propose a powerful spatial functional linear model which integrates the advantages of FLM in dealing with high dimensional data and SAR model in coping with spatial dependence.
- A simple estimation method is developed to obtain estimators
- Our simulation study demonstrate the consistency of the proposed estimators.
- A real dataset study demonstrates superiority of SFLM over FLM.

- We propose a powerful spatial functional linear model which integrates the advantages of FLM in dealing with high dimensional data and SAR model in coping with spatial dependence.
- A simple estimation method is developed to obtain estimators
- Our simulation study demonstrate the consistency of the proposed estimators.
- A real dataset study demonstrates superiority of SFLM over FLM.
- The SFLM with multiple functional predictors can be also obtained to deal with problem in practice.

- We propose a powerful spatial functional linear model which integrates the advantages of FLM in dealing with high dimensional data and SAR model in coping with spatial dependence.
- A simple estimation method is developed to obtain estimators
- Our simulation study demonstrate the consistency of the proposed estimators.
- A real dataset study demonstrates superiority of SFLM over FLM.
- The SFLM with multiple functional predictors can be also obtained to deal with problem in practice.
- We can also conduct functional variable selection in further research.

- Cai, T. and Hall, P. (2006). Prediction in functional linear regression. The Annals of Statistics, 34(5), 21592179.
- [2] Hall, P. and Horowitz, J. L. (2007). Methodology and convergence rates for functional linear regression. Annals of Statistics, 35(1), 7091.
- [3] Qu, X. and Lee, L. F. (2015). Estimating a spatial autoregressive model with an endogenous spatial weight matrix. Journal of Econometrics, 184(2), 209232.

Thank you very much !