

# Eigenvector-Based Semiparametric Filtering of Spatial Autocorrelation in Regression Models

Sebastian Juhl

University of Mannheim



March 17, 2021

- spatial autocorrelation can lead to inefficient or even biased and inconsistent parameter estimates in regression models
- social scientists primarily address this issue in one of two ways:
  - ① ignore it altogether and assume *iid* observations
  - ② apply (parametric) spatial regression models
- semiparametric filtering techniques can offer an attractive alternative
  - ease of estimation (standard OLS or ML estimators)
  - straightforward interpretation
  - accounts for spatial autocorrelation at different scales/resolutions
  - easily adaptable to GLMs

$$y = X\beta + e$$

$$y = X\beta + e$$

## Spatial Error DGP

$$y = X\beta + \underbrace{(I - \rho W)^{-1}\epsilon}_{e_{SEM}}$$

## Spatial-X DGP

$$y = X\beta + \underbrace{\rho W X + \epsilon}_{e_{SLX}}$$

## Spatial Lag DGP

$$y = X\beta + \underbrace{\rho W y + \epsilon}_{e_{SAR}}$$

$$y = (I - \rho W)^{-1}(X\beta + \epsilon)$$

$$y = X\beta + e$$

## Spatial Error DGP

$$y = X\beta + \underbrace{(I - \rho W)^{-1}\epsilon}_{e_{SEM}}$$

## Spatial-X DGP

$$y = X\beta + \underbrace{\rho W X + \epsilon}_{e_{SLX}}$$

## Spatial Lag DGP

$$y = X\beta + \underbrace{\rho W y + \epsilon}_{e_{SAR}}$$

$$y = (I - \rho W)^{-1}(X\beta + \epsilon)$$

$$(I - \rho W)^{-1} = (I + \rho W + \rho^2 W + \dots)$$

$$y = \left. \begin{array}{l} X\beta \\ + \rho W \epsilon \\ + \rho^2 W^2 \epsilon \\ + \dots + \epsilon \end{array} \right\} \epsilon_{SEM}$$

$$y = \left. \begin{array}{l} X\beta \\ + \rho W (X\beta + \epsilon) \\ + \rho^2 W^2 (X\beta + \epsilon) \\ + \dots + \epsilon \end{array} \right\} \epsilon_{SAR}$$

$$y = X\beta + e$$

## Spatial Error DGP

$$y = X\beta + \underbrace{(I - \rho W)^{-1}\epsilon}_{e_{SEM}}$$

## Spatial-X DGP

$$y = X\beta + \underbrace{\rho W X + \epsilon}_{e_{SLX}}$$

## Spatial Lag DGP

$$y = X\beta + \underbrace{\rho W y + \epsilon}_{e_{SAR}}$$

$$y = (I - \rho W)^{-1}(X\beta + \epsilon)$$

$$(I - \rho W)^{-1} = (I + \rho W + \rho^2 W + \dots)$$

$$y = \left. \begin{array}{l} X\beta \\ + \rho W \epsilon \\ + \rho^2 W^2 \epsilon \\ + \dots + \epsilon \end{array} \right\} \epsilon_{SEM}$$

$e_{SEM}$

$$\underbrace{\sum_{r=1}^{\infty} \rho^r W^r \epsilon}_{\text{spatial part}} + \underbrace{\epsilon}_{\text{noise}}$$

$$y = \left. \begin{array}{l} X\beta \\ + \rho W (X\beta + \epsilon) \\ + \rho^2 W^2 (X\beta + \epsilon) \\ + \dots + \epsilon \end{array} \right\} \epsilon_{SAR}$$

$e_{SAR}$

$$\underbrace{\sum_{r=1}^{\infty} \rho^r W^r (X\beta + \epsilon)}_{\text{spatial part}} + \underbrace{\epsilon}_{\text{noise}}$$

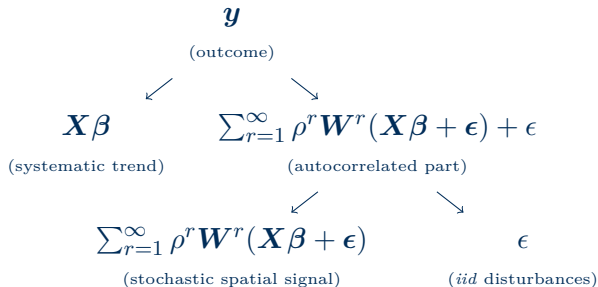
$e_{SLX}$

$$\underbrace{\rho W X}_{\text{spatial part}} + \underbrace{\epsilon}_{\text{noise}}$$

## Intuition:

ESF partitions the response variable into i) a systematic trend, ii) a stochastic spatial signal, and iii) *iid* disturbances

## Example: SAR DGP



ESF is based on the spectral decomposition of a centered (and symmetric/ symmetrized) connectivity matrix:  $\mathbf{W}$ :

$$\mathbf{M}\mathbf{W}\mathbf{M} = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^{-1} = \mathbf{E}\mathbf{\Lambda}\mathbf{E}'$$

- demeaning projector  $\mathbf{M} = (\mathbf{I} - \mathbf{1}\mathbf{1}'/n)$  – also ensures that all eigenvectors are orthogonal and uncorrelated
- $\mathbf{E}$  are all  $n$  eigenvectors
- $\mathbf{\Lambda}$  is a diagonal matrix of the corresponding eigenvalues  $\lambda$



ESF is based on the spectral decomposition of a centered (and symmetric/ symmetrized) connectivity matrix:  $W$ :

$$\textcolor{red}{MWM} = E\Lambda E^{-1} = E\Lambda E'$$

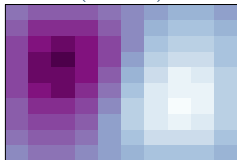
- demeaning projector  $M = (I - \mathbf{1}\mathbf{1}'/n)$  – also ensures that all eigenvectors are orthogonal and uncorrelated
- $E$  are all  $n$  eigenvectors
- $\Lambda$  is a diagonal matrix of the corresponding eigenvalues  $\lambda$

Direct relationship to the numerator of the global Moran coefficient:

$$MC(x) = \frac{n}{\mathbf{1}'W\mathbf{1}} \frac{x'\textcolor{red}{MWM}x}{x'Mx}$$

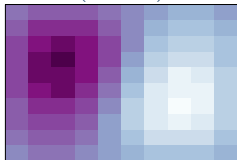
- eigenvectors  $\mathbf{E}$  depict distinct – and uncorrelated – synthetic map patterns
- corresponding eigenvalues indicate the level of SA

Eigenvector 1  
(MC = 1)

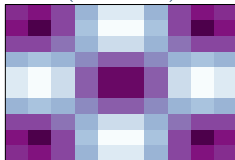


- eigenvectors  $\mathbf{E}$  depict distinct – and uncorrelated – synthetic map patterns
- corresponding eigenvalues indicate the level of SA

Eigenvector 1  
(MC = 1)

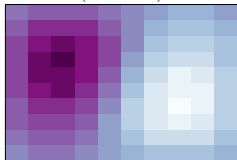


Eigenvector 10  
(MC = 0.73)

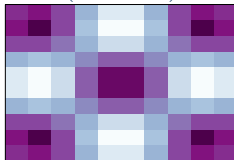


- eigenvectors  $\mathbf{E}$  depict distinct – and uncorrelated – synthetic map patterns
- corresponding eigenvalues indicate the level of SA

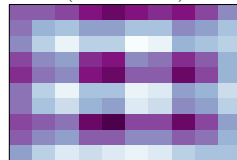
Eigenvector 1  
(MC = 1)



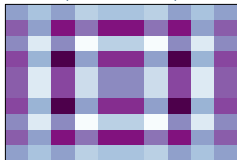
Eigenvector 10  
(MC = 0.73)



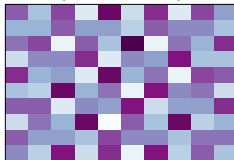
Eigenvector 20  
(MC = 0.454)



Eigenvector 44  
(MC = 0.066)

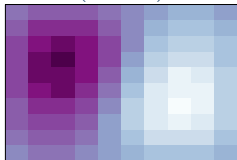


Eigenvector 80  
(MC = -0.454)

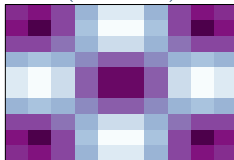


- eigenvectors  $E$  depict distinct – and uncorrelated – synthetic map patterns
- corresponding eigenvalues indicate the level of SA

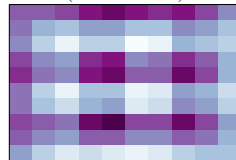
Eigenvector 1  
(MC = 1)



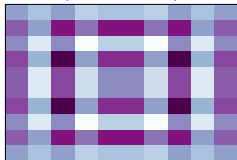
Eigenvector 10  
(MC = 0.73)



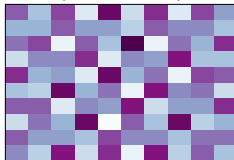
Eigenvector 20  
(MC = 0.454)



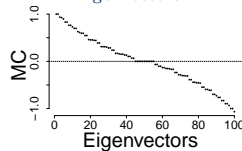
Eigenvector 44  
(MC = 0.066)



Eigenvector 80  
(MC = -0.454)



MCs of all  
Eigenvectors



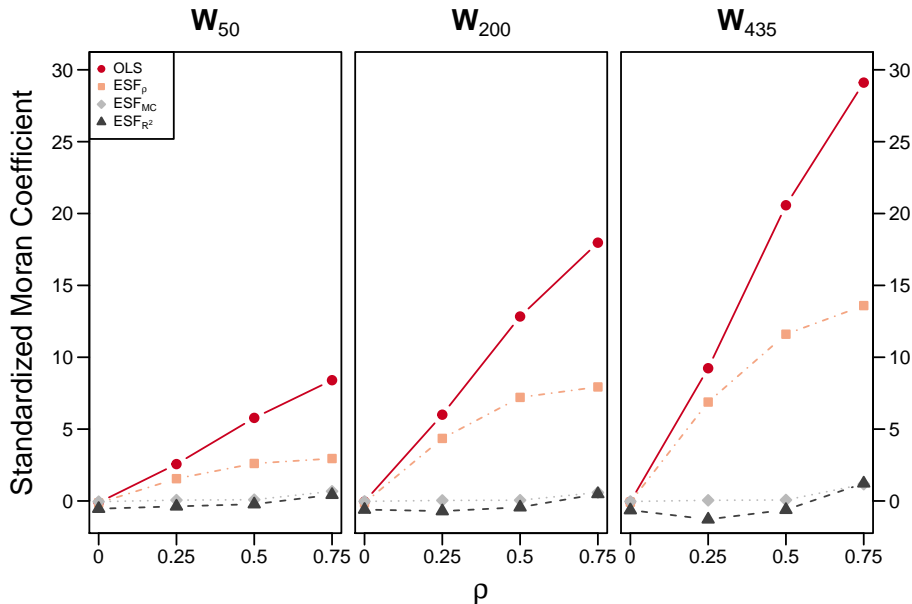
- $E$  depict all possible spatial patterns permitted by  $W$
- more complex patterns can be obtained by a linear combination of eigenvectors
- supervised or unsupervised stepwise regression techniques used to select the most relevant patterns  $E^*$

$$\begin{aligned} y &= X\beta + \overbrace{E\gamma}^{e_{OLS}} + \epsilon \\ &\approx X\beta + \underbrace{E^*\gamma}_{\text{filter}} + \underbrace{\epsilon}_{\text{noise}} \end{aligned}$$

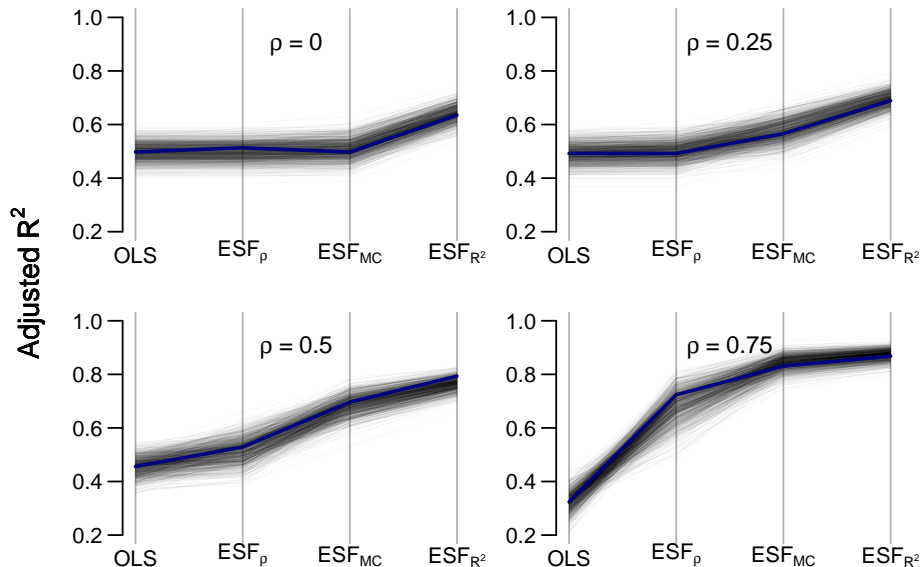
- different criteria can be used to select eigenvectors, e.g.:
  - maximization of model fit (e.g., AIC, BIC,  $R^2$ )
  - significance level of eigenvectors
  - minimization of residual SA

(Preliminary)

# Monte Carlo Evidence







- ESF offers numerous advantages for political scientists
- however, applicability depends on RQ (as always!)
- things to consider:

## Pros

- + ease of model estimation
- + straightforward interpretation of parameters
- + flexibility  
(no need to specify the spatial pattern in each variable)
- + generalizability  
(also applicable to GLMs – with some modifications)

## Cons

- “removes” indirect spillovers
- computationally demanding for large N
- over- or undercorrection of SA possible

- **spfilterR** package:  
CRAN: <https://CRAN.R-project.org/package=spfilterR>  
GitHub: <https://github.com/sjuhl/spfilterR>
- further projects & working papers
  - ❶ introductory paper on the **spfilterR** package (under review)
  - ❷ project on Moran eigenvector maps and spatial eigenfunction analysis

**Feedback & suggestions are highly appreciated!**

`sebastian.juhl@gess.uni-mannheim.de`

`www.sebastianjuhl.com`