

Local variable selection and parameter estimation for spatially varying coefficient models

Wesley Brooks

Department of Statistics
University of Wisconsin–Madison

October 31, 2013

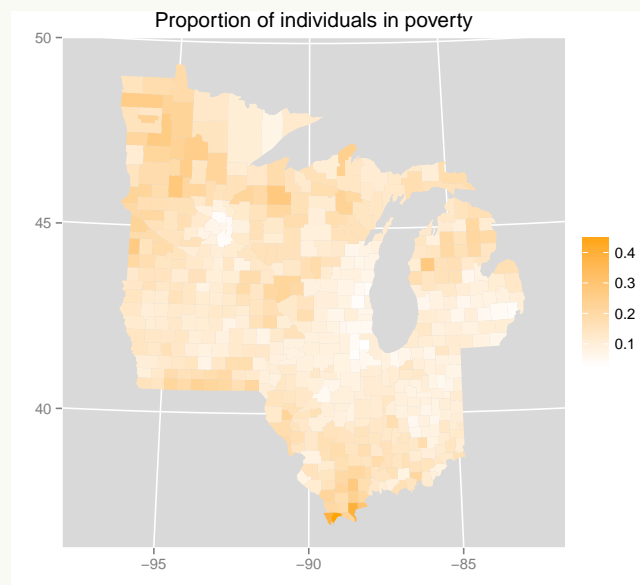
These slides were prepared for a practice version of my preliminary exam to advance to Ph.D candidacy in statistics at the University of Wisconsin–Madison.

Motivation

2

Motivation

Take a look at some data



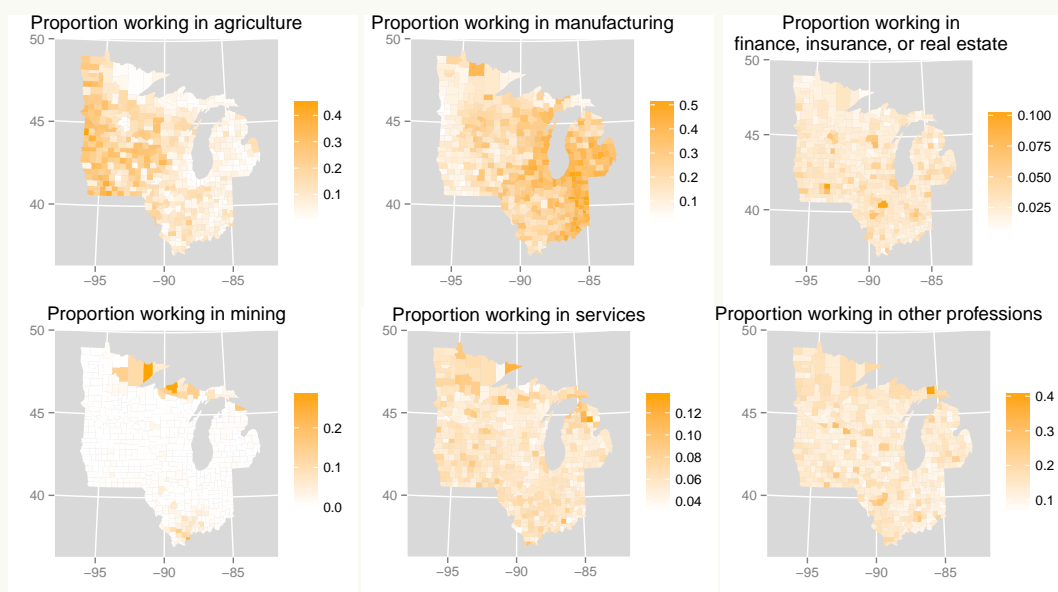
3

This is the county-level poverty rate from 1970, as well as the proportion of people who worked in manufacturing, agriculture, and services.

How is this data to be analyzed?

Motivation

Take a look at some data



This is the county-level poverty rate from 1970, as well as the proportion of people who worked in manufacturing, agriculture, and services.

How is this data to be analyzed?

Motivation

Sensible questions about the data

- ▶ Which of the economic-structure variables is associated with poverty rate?
- ▶ What are the sign and magnitude of that association?
- ▶ Is poverty rate associated with the same economic-structure variables across the entire region?
- ▶ Are the sign and magnitude of the associations constant across the region?

We're going to aim at answering these questions with the work I present today.

There are several other methods to answer at least some of these questions, which we'll cover next.

Introduction

Introduction

A review of existing methods

- ▶ Spatial regression
- ▶ Varying coefficient regression
 - Splines
 - Kernels
 - Wavelets
- ▶ Model selection via regularization

Behind the methodology that I'm discussing is a wide range of literature.

Introduction

Some definitions

- ▶ Univariate spatial response process $\{Y(\boldsymbol{s}) : \boldsymbol{s} \in \mathcal{D}\}$
- ▶ Multivariate spatial covariate process $\{\boldsymbol{X}(\boldsymbol{s}) : \boldsymbol{s} \in \mathcal{D}\}$
- ▶ n = number of observations
- ▶ p = number of covariates
- ▶ Location (2-dimensional) \boldsymbol{s}
- ▶ Spatial domain \mathcal{D}

We'll use these variables throughout.

Introduction

Further definitions

► Geostatistical data:

- Observations are made at sampling locations s_i for $i = 1, \dots, n$
- E.g. elevation, temperature

► Areal data:

- Domain is partitioned into n regions $\{D_1, \dots, D_n\}$
- The regions do not overlap, and they divide the domain completely: $\mathcal{D} = \bigcup_{i=1}^n D_i$
- Sampling locations s_i for $i = 1, \dots, n$ are the centroids of the regions
- E.g. poverty rate, population, spatial mean temperature

The method I'm describing applies to geostatistical data, or to areal data when the observations are assumed to be located at the centroid.

The poverty data example is areal data, the simulation study is based on simulated geostatistical data.

Introduction

Varying coefficients regression (Hastie and Tibshirani, 1993)

$$Y(\boldsymbol{s}) = \boldsymbol{X}(\boldsymbol{s})'\boldsymbol{\beta}(\boldsymbol{s}) + \varepsilon(\boldsymbol{s})$$

- ▶ Assume an effect modifying variable S
- ▶ Coefficients are functions of S

note

Introduction

Spatial regression

- ▶ The typical spatial regression (Cressie, 1993)

$$\begin{aligned} Y(\mathbf{s}) &= \mathbf{X}(\mathbf{s})'\boldsymbol{\beta} + W(\mathbf{s}) + \varepsilon(\mathbf{s}) \\ \text{cov}(W(\mathbf{s}), W(\mathbf{t})) &= \Gamma(\delta(\mathbf{s}, \mathbf{t})) \\ \delta(\mathbf{s}, \mathbf{t}) &= \sqrt{\|\mathbf{t} - \mathbf{s}\|_2} \\ \text{E.g. } \Gamma(\delta(\mathbf{s}, \mathbf{t})) &= \exp\{-\phi^{-1}\delta(\mathbf{s}, \mathbf{t})\} \end{aligned} \tag{1}$$

- ▶ $W(\mathbf{s})$ is a spatial random effect that accounts for autocorrelation in the response variable
- ▶ The coefficients $\boldsymbol{\beta}$ are constant
- ▶ Relies on *a priori* global variable selection

11

This is the form of the usual spatial regression from e.g. Cressie (1993).

The spatial random effect W captures autocorrelation of the response, while the white noise is iid error

The Gamma function is a Matern-class covariance function, such as the exponential covariance function (listed here)

Introduction

Spatially varying coefficient process (Gelfand et al., 2003)

- ▶ Making model more flexible: coefficients in a spatial regression model can be allowed to vary

$$Y(\mathbf{s}) = \mathbf{X}(\mathbf{s})'\boldsymbol{\beta}(\mathbf{s}) + \varepsilon(\mathbf{s})$$

- ▶ The spatial random effect has been incorporated into the spatially varying intercept
- ▶ $\{\beta_1(\mathbf{s}) : \mathbf{s} \in \mathcal{D}\}, \dots, \{\beta_p(\mathbf{s}) : \mathbf{s} \in \mathcal{D}\}$ are stationary spatial processes with Matérn covariance functions
- ▶ Still relies on *a priori* global variable selection

note

Introduction

Spline-based VCR models (Wood, 2006)

- ▶ Splines are a way to parameterize smooth functions
- ▶ Splines can be incorporated into a generalized additive model (GAM):
 - $E\{Y(t)\} = f\{X_1(t)\} + \cdots + f\{X_p(t)\}$
- ▶ It is possible to parameterize a VCR model with splines for the coefficient functions:
 - $E\{Y(t)\} = \beta_1(t)X_1(t) + \cdots + \beta_p(t)X_p(t)$

13

note

Introduction

Global selection in spline-based VCR models (Wang, H. Li, and Huang, 2008; Antoniadis, Gijbels, and Verhasselt, 2012)

Regularization methods for global variable selection in VCR models:

- ▶ The integral of a function squared (e.g. $\int \{f(t)\}^2 dt$) is zero if and only if the function is zero everywhere.
- ▶ Use regularization (maximize the likelihood plus a penalty) to encourage coefficient functions to be zero
- ▶ SCAD penalty (Fan and R. Li, 2001) on the integral of the square of the coefficient function
- ▶ Non-negative garrote penalty (Breiman, 1995) on the integral of the square of the coefficient function

14

These selection methods are all global - that is, they select variables for the entire domain simultaneously

Introduction

Existing approaches: wavelet methods for VCR models

Wavelet methods involve decomposing a function into local frequency components. Wavelet methods for VCR models include using Bayesian variable selection or the Lasso to estimate which local frequency components have nonzero coefficients (Shang, 2011; J. Zhang and Clayton, 2011).

These methods achieve sparsity in the local frequency components but not in the local covariates, and so are not suitable for local variable selection.

15

note

Geographically weighted regression

16

Geographically weighted regression

Existing approaches: geographically weighted regression

When the effect modifying variable s refers to spatial location, the method of local regression is called geographically weighted regression (GWR) (Brundson, S. Fotheringham, and Martin Charlton, 1998; A. Fotheringham, Brunsdon, and M. Charlton, 2002)

17

Maximizing the local likelihood for a model of Gaussian data with iid errors can be done by weighted least squares.

Geographically weighted regression

Existing approaches: Local regression

Local regression uses a kernel function at each sampling location to weight observations based on their distance from the sampling location. An example is the bisquare kernel:

$$w_{ii'} = \begin{cases} \left[1 - (\phi^{-1}\delta_{ii'})^2\right]^2 & \text{if } \delta_{ii'} < \phi, \\ 0 & \text{if } \delta_{ii'} \geq \phi. \end{cases} \quad (2)$$

Where ϕ is a bandwidth parameter.

Given the weights, a local model is fit at each sampling location using the local likelihood (Loader, 1999)

note

Geographically weighted regression

Existing approaches: Local likelihood

Calibrate the model by doing the following at each sampling location:

- ▶ Weight each observation's likelihood
- ▶ Weights are given by the kernel

$$L = \prod_{i'=1}^n (L_{i'})^{w_{ii'}}$$
$$\ell = \sum_{i'=1}^n w_{ii'} \left\{ \log \sigma_i^2 + \sigma_i^{-2} (y_{i'} - \mathbf{x}_{i'}' \boldsymbol{\beta}_i)^2 \right\}$$

Where $\boldsymbol{\beta}_i = \boldsymbol{\beta}(\mathbf{s}_i)$.

Maximizing the local likelihood for a model of Gaussian data with iid errors can be done by weighted least squares.

Geographically weighted regression

Existing approaches: bandwidth estimation for GWR

- ▶ Smaller bandwidth: less bias, more flexible coefficient surface
- ▶ Large bandwidth: less variance, less flexible coefficient surface
- ▶ Estimate the degrees of freedom used in estimating the coefficient surface (Hurvich, Simonoff, and Tsai, 1998):
 - $\hat{y} = Hy$
 - $\nu = \text{tr}(H)$
- ▶ Then the corrected AIC for bandwidth selection is:
- ▶ $\text{AIC}_c = 2n \log \sigma + n \left\{ \frac{n+\nu}{n-2-\nu} \right\}$

Maximizing the local likelihood for a model of Gaussian data with iid errors can be done by weighted least squares.

Geographically weighted regression

Existing approaches: geographically weighted Lasso

Within a GWR model, using the Lasso (Tibshirani, 1996) for local variable selection is called the geographically weighted Lasso (GWL) (Wheeler, 2009).

- ▶ The GWL requires estimating a Lasso tuning parameter for each local model
- ▶ Wheeler, 2009 estimates the local Lasso tuning parameter at location s_i by minimizing a jackknife criterion: $|y_i - \hat{y}_i|$
- ▶ The jackknife criterion can only be calculated where data are observed, making it impossible to use the GWL to impute missing data or to estimate the value of the coefficient surface at new locations
- ▶ Also, the Lasso is known to be biased in variable selection and suboptimal for coefficient estimation

GWL does local variable selection

Local variable selection and parameter
estimation

Local variable selection and parameter estimation

Geographically weighted elastic net (GWEN)

- ▶ Local variable selection in a GWR model using the adaptive elastic net (AEN) (Zou and H. Zhang, 2009)
- ▶ Under suitable conditions, the AEN has an oracle property for selection

$$\begin{aligned}\mathcal{S}(\boldsymbol{\beta}_i) &= -2\ell_i(\boldsymbol{\beta}_i) + \mathcal{J}_2(\boldsymbol{\beta}_i) \\ &= \sum_{i'=1}^n w_{ii'} \left\{ \log \sigma_i^2 + (\sigma_i^2)^{-1} (y_{i'} - \mathbf{x}_{i'}' \boldsymbol{\beta}_i)^2 \right\} \\ &\quad + \alpha_i \lambda_i \sum_{j=1}^p |\beta_{ij}| / \gamma_{ij} \\ &\quad + (1 - \alpha_i) \lambda_i^* \sum_{j=1}^p (\beta_{ij} / \gamma_{ij})^2\end{aligned}$$

23

The adaptive weights $\boldsymbol{\gamma}_i = (\gamma_{i1}, \dots, \gamma_{ip})'$ are defined in the same way as for the AL, and the elastic net parameter $\alpha_i \in [0, 1]$ controls the balance between ℓ_1 penalty $\lambda_i^* \sum_{j=1}^p |\beta_{ij}| / \gamma_{ij}$ and ℓ_2 penalty $\lambda_i^* \sum_{j=1}^p (\beta_{ij} / \gamma_{ij})^2$.

Local variable selection and parameter estimation

Geographically weighted elastic net (GWEN)

where $\sum_{i'=1}^n w_{ii'} (y_{i'} - \mathbf{x}_{i'}' \boldsymbol{\beta}_i)^2$ is the weighted sum of squares minimized by traditional GWR, and $\mathcal{J}_2(\boldsymbol{\beta}_i) = \alpha_i \lambda_i^* \sum_{j=1}^p |\beta_{ij}| / \gamma_{ij} + (1 - \alpha_i) \lambda_i^* \sum_{j=1}^p (\beta_{ij} / \gamma_{ij})^2$ is the AEN penalty.

24

The adaptive weights $\boldsymbol{\gamma}_i = (\gamma_{i1}, \dots, \gamma_{ip})'$ are defined in the same way as for the AL, and the elastic net parameter $\alpha_i \in [0, 1]$ controls the balance between ℓ_1 penalty $\lambda_i^* \sum_{j=1}^p |\beta_{ij}| / \gamma_{ij}$ and ℓ_2 penalty $\lambda_i^* \sum_{j=1}^p (\beta_{ij} / \gamma_{ij})^2$.

Local variable selection and parameter estimation

Geographically weighted elastic net (GWEN)

It is necessary to estimate an AEN tuning parameter for each local model. Using the local BIC allows fitting a local model at any location within the domain

$$\begin{aligned} \text{BIC}_{\text{loc},i} &= -2 \sum_{i'=1}^n \ell_{ii'} + \log \left(\sum_{i'=1}^n w_{ii'} \right) \text{df}_i \\ &= -2 \sum_{i'=1}^n \log \left[(2\pi \hat{\sigma}_i^2)^{-1/2} \exp \left\{ -\frac{1}{2} \hat{\sigma}_i^{-2} \left(y_{i'} - \mathbf{x}_{i'}' \hat{\boldsymbol{\beta}}_{i'} \right)^2 \right\} \right]^w \\ &\quad + \log \left(\sum_{i'=1}^n w_{ii'} \right) \text{df}_i \end{aligned} \tag{3}$$

25

The adaptive weights $\boldsymbol{\gamma}_i = (\gamma_{i1}, \dots, \gamma_{ip})'$ are defined in the same way as for the AL, and the elastic net parameter $\alpha_i \in [0, 1]$ controls the balance between ℓ_1 penalty $\lambda_i^* \sum_{j=1}^p |\beta_{ij}|/\gamma_{ij}$ and ℓ_2 penalty $\lambda_i^* \sum_{j=1}^p (\beta_{ij}/\gamma_{ij})^2$.

Local variable selection and parameter estimation

Geographically weighted elastic net (GWEN)

$$\begin{aligned} &= \sum_{i'=1}^n w_{ii'} \left\{ \log(2\pi) + \log \hat{\sigma}_i^2 + \hat{\sigma}_i^{-2} \left(y_{i'} - \mathbf{x}_{i'}' \hat{\boldsymbol{\beta}}_{i'} \right)^2 \right\} \\ &+ \log \left(\sum_{i'=1}^n w_{ii'} \right) \text{df}_i \end{aligned}$$

26

We treat the sum of the weights around the sampling location as the number of observations for the local BIC.

Simulation study

27

Simulation study

Simulating covariates

Five covariates $\tilde{X}_1, \dots, \tilde{X}_5$ were simulated by Gaussian random fields on the domain $[0, 1] \times [0, 1]$ on a 30×30 grid of sampling locations:

$$\begin{aligned} \tilde{X}_j &\sim N(0, \Sigma) \text{ for } j = 1, \dots, 5 \\ \{\Sigma\}_{i,i'} &= \exp\{-\tau^{-1}\delta_{ii'}\} \text{ for } i, i' = 1, \dots, n \end{aligned}$$

Where the covariates were simulated with colinearity, the colinearity was induced by multiplying the design matrix by the square root of the colinearity matrix:

$$\begin{aligned} \text{diag}(\Omega_{5 \times 5}) &= 1 \\ \text{off-diag}(\Omega_{5 \times 5}) &= \rho \\ X &= \tilde{X}R \end{aligned} \tag{4}$$

Where $\Omega_{5 \times 5} = R'R$ is the Cholesky decomposition.

28

note

Simulation study

Simulating the response

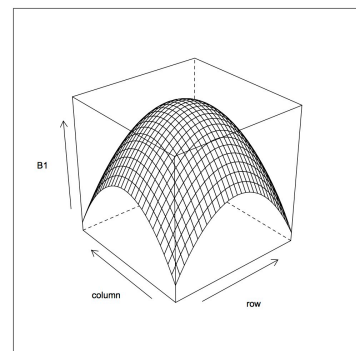
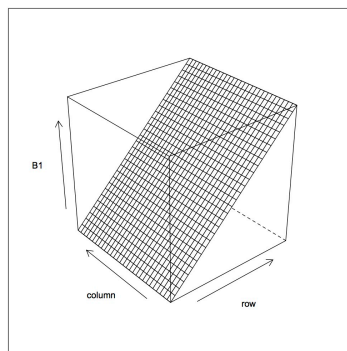
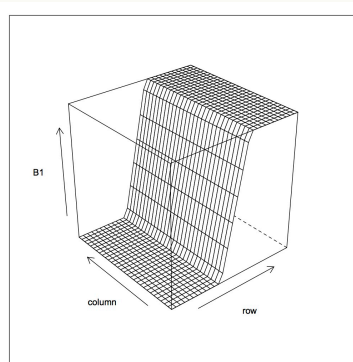
- ▶ $Y(\mathbf{s}) = X(\mathbf{s})'\boldsymbol{\beta}(\mathbf{s}) = \sum_{j=1}^5 \beta_j(\mathbf{s})X_j(\mathbf{s}) + \varepsilon(\mathbf{s})$
- ▶ $\varepsilon \sim iid\ N(0, \sigma^2)$
- ▶ $\beta_1(\mathbf{s})$, the coefficient function for X_1 , is nonzero in part of the domain.
- ▶ Coefficients for X_2, \dots, X_5 are zero everywhere

note

Simulation study

Coefficient functions

Call these functions step, gradient, and parabola:



note

Simulation study

Simulation settings

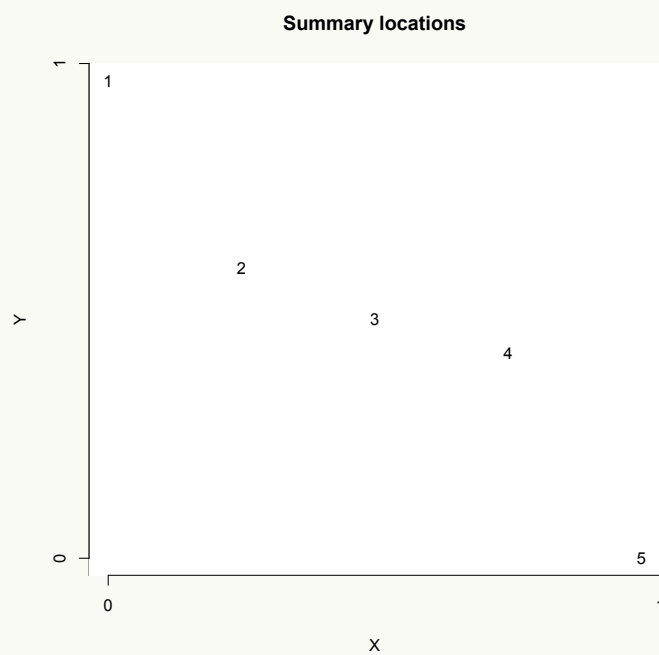
Setting	function	ρ	σ^2
1	step	0	0.25
2	step	0	1
3	step	0.5	0.25
4	step	0.5	1
5	gradient	0	0.25
6	gradient	0	1
7	gradient	0.5	0.25
8	gradient	0.5	1
9	parabola	0	0.25
10	parabola	0	1
11	parabola	0.5	0.25
12	parabola	0.5	1

Table : Simulation parameters for each setting.

note

Simulation results

Selection



note

Simulation results

Selection

location	step				gradient				parabola			
	GWEN		GWAL		GWEN		GWAL		GWEN		GWAL	
	β_1	$\beta_2 - \beta_5$	β_1	$\beta_2 - \beta_5$	β_1	$\beta_2 - \beta_5$	β_1	$\beta_2 - \beta_5$	β_1	$\beta_2 - \beta_5$	β_1	$\beta_2 - \beta_5$
1	0.99	0.00	0.99	0.00	1.00	0.00	1.00	0.00	0.36	0.00	0.38	0.00
	0.99	0.02	0.99	0.02	1.00	0.01	1.00	0.01	0.71	0.02	0.70	0.02
	0.99	0.00	1.00	0.00	1.00	0.00	1.00	0.00	0.28	0.00	0.33	0.00
	0.96	0.05	0.91	0.04	0.99	0.03	0.99	0.01	0.56	0.02	0.55	0.02
2	1.00	0.00	1.00	0.00	1.00	0.01	1.00	0.00	1.00	0.00	1.00	0.00
	1.00	0.03	1.00	0.03	1.00	0.02	1.00	0.02	1.00	0.02	0.99	0.01
	1.00	0.01	1.00	0.00	1.00	0.01	1.00	0.01	1.00	0.00	1.00	0.00
	0.99	0.05	0.97	0.04	1.00	0.02	0.99	0.01	0.98	0.02	0.97	0.01
3	0.91	0.01	0.91	0.00	1.00	0.01	1.00	0.01	1.00	0.00	1.00	0.00
	0.96	0.05	0.96	0.05	1.00	0.01	1.00	0.01	1.00	0.00	1.00	0.00
	0.92	0.05	0.95	0.02	1.00	0.02	1.00	0.01	1.00	0.00	1.00	0.00
	0.92	0.08	0.87	0.05	1.00	0.02	0.98	0.02	0.99	0.01	0.99	0.01
4	0.48	0.01	0.43	0.01	1.00	0.01	1.00	0.01	1.00	0.00	1.00	0.00
	0.72	0.04	0.78	0.03	1.00	0.01	1.00	0.01	1.00	0.00	1.00	0.00
	0.49	0.02	0.46	0.02	1.00	0.02	1.00	0.00	1.00	0.00	1.00	0.00
	0.60	0.05	0.56	0.04	1.00	0.03	0.98	0.02	1.00	0.01	0.98	0.02
5	0.00	0.00	0.00	0.00	0.83	0.00	0.82	0.00	0.32	0.00	0.32	0.00
	0.03	0.01	0.02	0.00	0.70	0.00	0.66	0.00	0.68	0.02	0.73	0.02
	0.00	0.00	0.00	0.00	0.87	0.01	0.87	0.00	0.37	0.00	0.42	0.00
	0.06	0.02	0.01	0.02	0.61	0.01	0.62	0.02	0.61	0.04	0.58	0.03

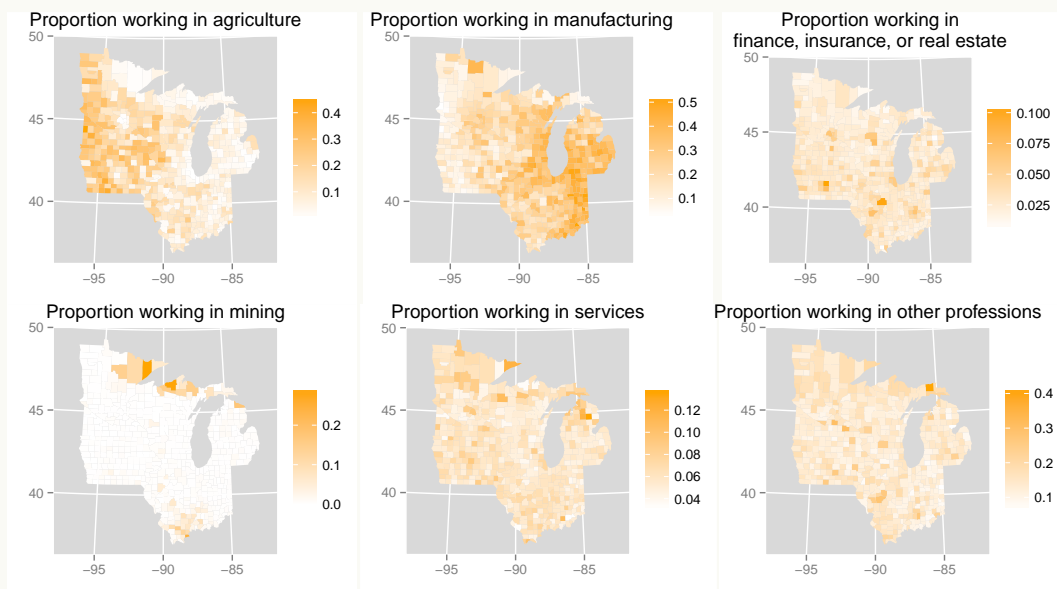
Table : Selection frequency for the indicated variables.

note

Data example: poverty rate in the upper
midwest

Data example: poverty rate in the upper midwest

Revisiting the introductory example



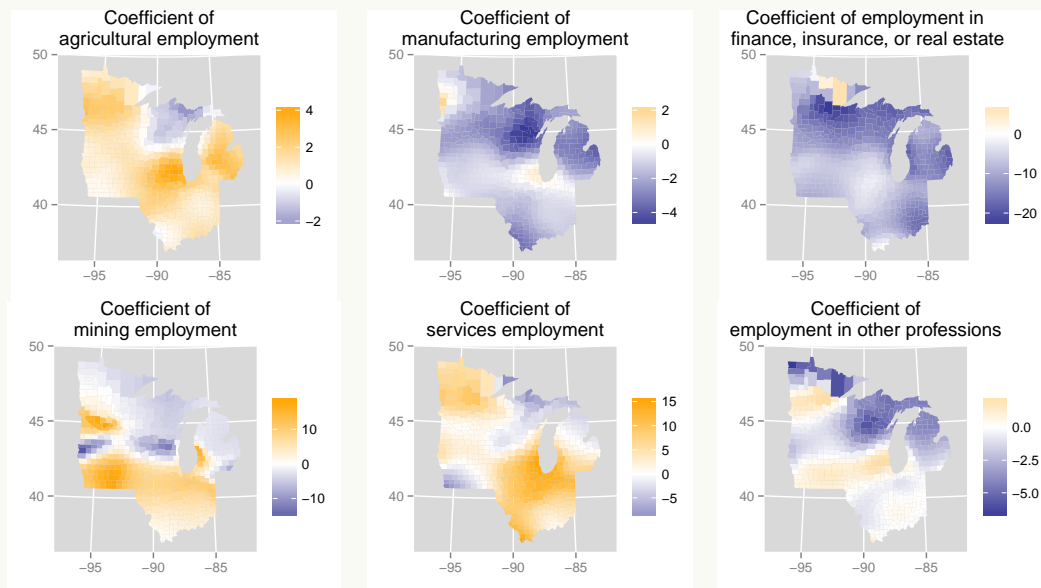
35

This is the county-level poverty rate from 1970, as well as the proportion of people who worked in manufacturing, agriculture, and services.

How is this data to be analyzed?

Data example: poverty rate in the upper midwest

Results from traditional GWR



note

Data example: poverty rate in the upper midwest

Data description

- ▶ Response: logit-transformed poverty rate in the Upper Midwest states of the U.S.
 - Minnesota, Iowa, Wisconsin, Illinois, Indiana, Michigan
- ▶ Covariates: employment structure (raw proportion employed in:)
 - agriculture
 - finance, insurance, and real estate
 - manufacturing
 - mining
 - services
 - other professions
- ▶ Data source: U.S. Census Bureau's decennial census of 1970

37

note

Data example: poverty rate in the upper midwest

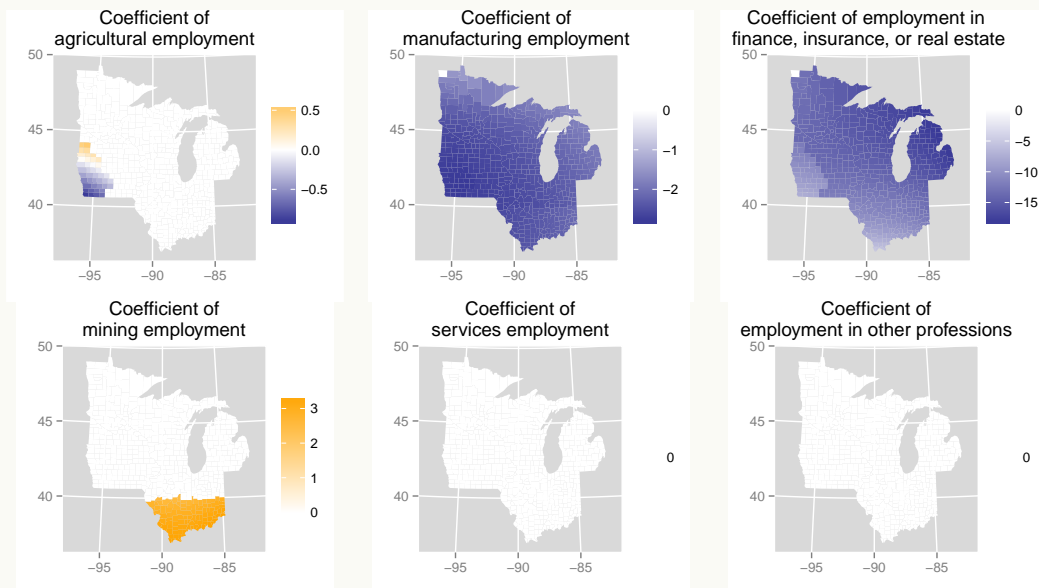
Data description

- ▶ Data aggregated to the county level
 - counties are areal units
- ▶ county centroid treated as sampling location

note

Data example: poverty rate in the upper midwest

Results from GWEN



note

Future work

40

Future work

- ▶ Apply the GWEN to data with non-Gaussian response variable
- ▶ Incorporate spatial autocorrelation in the model (simulated errors were iid)

41

note

Acknowledgements