

Local variable selection and coefficient estimation for spatially-varying-coefficients models in the context of geographically-weighted regression using regularization

Wesley Brooks

1. Introduction

Varying-coefficients regression (VCR) (Hastie and Tibshirani, 1993) is a technique used to model non-stationary regression processes. Whereas the coefficients in ordinary least squares (OLS) regression are scalar constants, the coefficients in VCR are functions - often *smooth* functions - of some effect-modifying parameter. Methods for estimating the coefficient functions of VCR are typically divided between spline-based (Wood, 2006) and kernel-based (Hastie and Loader, 1993; Loader, 1999) methods.

Geostatistical data is the name for observations of a continuous spatial process that are made at discrete locations. Clustering is a typical form of non-stationarity in spatial data. One way to analyze clustered data is to assume a constant underlying mean with random deviations that are clustered in space. This is the structure of an autoregressive model (). If the underlying mean is not constant but is given by a regression function with constant coefficients, then a conditionally autoregressive (CAR) model () is used to simultaneously estimate regression coefficients and clustering of the residuals.

Spatially-clustered residuals may indicate that the random error component arises in a spatially-clustered way but it may also indicate oversmoothing, where the small net bias of a global estimate masks larger but offsetting local bias. In this case, smoothing methods that model the mean

response as a function of the covariates can reduce or eliminate clustering of the residuals. Where the underlying mean is a regression function but the coefficients are not constants, a spatial VCR model is appropriate. Both spline-based () and kernel-based (Fotheringham et al., 2002) methods are available for estimating the coefficient functions.

Geographically weighted regression (GWR) (Fotheringham et al., 2002) is a kernel-based method of estimating the coefficients of a VCR model for spatial data. GWR uses kernel-weighted regression with weights based on the distance between observation locations. The presentation of GWR in Fotheringham et al. (2002) follows the development of local likelihood in Loader (1999).

GWR can be thought of as a kernel smoother for regression coefficients, and hence GWR coefficient estimates are likely to exhibit bias near the boundary of the region being modeled (Hastie and Loader, 1993). Modeling the coefficient surface as locally linear rather than locally constant (by including coefficient-by-location interactions) can reduce this boundary-effect bias (Hastie and Loader, 1993). Adding these interactions to the GWR model is analogous to a transition from kernel smoothing to local regression, and was introduced in Wang et al. (2008).

There is interest among practitioners of GWR not only in estimating the spatially-varying coefficient surface, but in doing variable selection to estimate which covariates are important predictors of the output variable, and in what regions they are important (citations?).

Some recent research has focused on variable selection in varying-coefficients models. In the context of varying-coefficients regression models, global variable selection (in which one compares the hypothesis that the coefficient on a given variable is zero everywhere against the hypothesis that the coefficient is nonzero somewhere) is distinguished from local variable selection (in which one compares the hypothesis that the coefficient on a given variable is zero at a given location against the hypothesis that the coefficient at that location is nonzero). Global variable selection for models where the varying coefficients are estimated using splines is addressed in Fan and Zhang (1999) for response variables that belong to an exponential-family distribution (as in the generalized linear

model), and in Wang et al. (2008) for models with repeated measurements. Antoniadis et al. (2012) estimates the coefficient functions with P-splines, and then uses the nonnegative garrote of Breiman (1995) to do local variable selection by selecting P-spline bases.

Here we discuss a method of local variable selection in GWR models using the adaptive LASSO of Zou (2006). The idea first appears in the literature as the geographically-weighted LASSO (GWL) of Wheeler (2009), which uses a jackknife criterion for selection of the LASSO tuning parameters. Because the jackknife criterion can only be computed at locations where the response variable is observed, the GWL cannot be used for imputation of missing data nor for interpolation between observation locations. We avoid this limitation of the GWL by using a penalized-likelihood criterion to select the LASSO tuning parameters (specifically the AIC, but in principle one could use the BIC, *et cetera*). The AIC allows us to easily adapt our method to the setting of a generalized linear model. The local AIC presented here is based on an *ad hoc* calculation of the sample size and degrees of freedom for estimating the spatially-varying coefficient surfaces.

2. Geographically-weighted regression

2.1. Model

Consider n data observations, made at sampling locations $\mathbf{s}_1, \dots, \mathbf{s}_n$ in a spatial domain $D \subset \mathbb{R}^2$. For $i = 1, \dots, n$, let $y(\mathbf{s}_i)$ and $\mathbf{x}(\mathbf{s}_i)$ denote the univariate response variable, and a $(p + 1)$ -variate vector of covariates measured at location \mathbf{s}_i , respectively. At each location \mathbf{s}_i , assume that the outcome is related to the covariates by a linear model where the coefficients $\boldsymbol{\beta}_i(\mathbf{s}_i)$ may be spatially-varying and $\varepsilon(\mathbf{s}_i)$ is random noise at location \mathbf{s}_i .

$$y(\mathbf{s}_i) = \mathbf{x}'(\mathbf{s}_i)\boldsymbol{\beta}_i(\mathbf{s}_i) + \varepsilon(\mathbf{s}_i) \tag{1}$$

Further assume that the error term $\varepsilon(\mathbf{s}_i)$ is normally distributed with zero mean and a possibly

spatially-varying variance $\sigma^2(\mathbf{s}_i)$

$$\varepsilon(\mathbf{s}_i) \sim \mathcal{N}(0, \sigma^2(\mathbf{s}_i)) \quad (2)$$

In order to simplify the notation, let subscripts denote the values of data or parameters at the locations where data is observed. Thus, $\mathbf{x}(\mathbf{s}_i) \equiv \mathbf{x}_i \equiv (1, x_{i1}, \dots, x_{ip})'$, $\boldsymbol{\beta}(\mathbf{s}_i) \equiv \boldsymbol{\beta}_i \equiv (\beta_{i0}, \beta_{i1}, \dots, \beta_{ip})'$, $y(\mathbf{s}_i) \equiv y_i$, and $\sigma^2(\mathbf{s}_i) \equiv \sigma_i^2$. Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$ and $\mathbf{Y} = (y_1, \dots, y_n)'$. Equations (1) and (2) can now be rewritten as

$$y_i = \mathbf{x}_i' \boldsymbol{\beta}_i + \epsilon_i \text{ and } \epsilon_i \sim \mathcal{N}(0, \sigma_i^2) \quad (3)$$

Assume that, given the design matrix \mathbf{X} , observations of the response variable at different locations are statistically independent of each other. Then the total log-likelihood of the observed data is the sum of the log-likelihood of each individual observation.

$$\ell(\boldsymbol{\beta}) = -1/2 \sum_{i=1}^n \left\{ \log(2\pi\sigma_i^2) + \sigma_i^{-2} (y_i - \mathbf{x}_i' \boldsymbol{\beta}_i)^2 \right\} \quad (4)$$

With n observations and $n \times (p+1)$ free parameters, the model is not identifiable so it is not possible to directly maximize the total likelihood. One way to effectively reduce the number of parameters is to assume that the spatially-varying coefficients $\boldsymbol{\beta}(\mathbf{s})$ are smoothly varying, and use a kernel smoother to make pointwise estimates of the coefficients by maximizing the local likelihood. In the setting of spatial data and with the kernel smoother based on the physical distance between observation locations, this is ordinary GWR.

2.2. Estimation

Geographically-weighted regression estimates the value of the coefficient surface $\boldsymbol{\beta}(\mathbf{s})$ at each location \mathbf{s}_i . First calculate the euclidean distance $\delta_{ii'} \equiv \delta(\mathbf{s}_i, \mathbf{s}_{i'}) \equiv \|\mathbf{s}_i - \mathbf{s}_{i'}\|_2$ between locations \mathbf{s}_i and $\mathbf{s}_{i'}$ for all i, i' . A bisquare kernel is used to generate spatial weights based on the euclidean

distances and a bandwidth ϕ :

$$w_{ii'} = \begin{cases} \left[1 - (\phi^{-1}\delta_{ii'})^2\right]^2 & \text{if } \delta_{ii'} < \phi \\ 0 & \text{if } \delta_{ii'} \geq \phi \end{cases} \quad (5)$$

For the purpose of estimation, define the local likelihood at each location (Fotheringham et al., 2002):

$$\mathcal{L}_i(\beta_i) = \prod_{i'=1}^n \left\{ (2\pi\sigma_i^2)^{-1/2} \exp \left[-\frac{1}{2}\sigma_i^{-2} (y_{i'} - \mathbf{x}_{i'}'\beta_i)^2 \right] \right\}^{w_{ii'}} \quad (6)$$

Thus, the local log-likelihood function is:

$$\ell_i(\beta_i) \propto -1/2 \sum_{i'=1}^n w_{ii'} \left\{ \log \sigma_i^2 + \sigma_i^{-2} (y_{i'} - \mathbf{x}_{i'}'\beta_i)^2 \right\} \quad (7)$$

From which it is apparent that the GWR coefficient estimates $\hat{\beta}_{i,\text{GWR}}$, which maximize the local likelihood at location \mathbf{s}_i , can be calculated using weighted least squares. Letting the diagonal weight matrix \mathbf{W}_i be:

$$\mathbf{W}_i = \text{diag} \{w_{ii'}\}_{i'=1}^n \quad (8)$$

We have:

$$\hat{\beta}_{i,\text{GWR}} = (\mathbf{X}'\mathbf{W}_i\mathbf{X})^{-1} \mathbf{X}'\mathbf{W}_i\mathbf{Y} \quad (9)$$

And $\hat{\sigma}_i$, which maximizes (7), is:

$$\hat{\sigma}_i = (\mathbf{1}_n' \mathbf{w}_i)^{-1} \mathbf{w}_i' \left(\mathbf{Y} - \mathbf{X} (\mathbf{X}'\mathbf{W}_i\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} \right) \quad (10)$$

3. Model selection

3.1. Variable selection

Traditional GWR relies on *a priori* model selection to decide which variables should be included in the model. In the context of ordinary least squares regression, regularization methods such as the adaptive LASSO (Zou, 2006) have been shown to have appealing properties for automating variable selection, sometimes including the “oracle” property of asymptotically selecting exactly the correct variables for inclusion in a regression model.

The adaptive LASSO is applied to GWR by first multiplying the design matrix \mathbf{X} by $\mathbf{W}_i^{1/2}$, the diagonal matrix of geographic weights centered at s_i . Since some of the weights $w_{ii'}$ may be zero, the matrix $\mathbf{W}_i^{1/2}\mathbf{X}$ is not of full rank. The matrices \mathbf{Y}_i^* , \mathbf{X}_i^* , and \mathbf{W}_i^* are formed by dropping the rows of \mathbf{X} and \mathbf{W}_i that correspond to observations with zero weight in the regression model at location s_i . Now, letting $\mathbf{U}_i^* = \mathbf{W}_i^{*1/2}\mathbf{X}_i^*$ and $\mathbf{V}_i^* = \mathbf{W}_i^{*1/2}\mathbf{Y}_i^*$, we seek the coefficients β_i of the regression model:

$$\mathbf{V}_i^* = \mathbf{U}_i^* \beta_i + \epsilon \quad (11)$$

To apply the adaptive LASSO for estimating these regression coefficients, each column of \mathbf{U}_i^* is centered around zero and rescaled to have an L_2 -norm of one. Let $\tilde{\mathbf{U}}_i^*$ be the centered-and-scaled version of \mathbf{U}_i^* . Adaptive weights are calculated using the OLS regression coefficients γ_i^* via ordinary least squares (OLS):

$$\gamma_i^* = \left(\tilde{\mathbf{U}}_i^{*'} \tilde{\mathbf{U}}_i^* \right)^{-1} \tilde{\mathbf{U}}_i^{*'} \mathbf{V}_i^* \quad (12)$$

Now a final scaling step is done: for $j = 1, \dots, p$, the j th column of $\tilde{\mathbf{U}}_i^*$ is multiplied by $(\gamma_i^*)_j$, the corresponding coefficient from (12). Call this rescaled matrix $\check{\mathbf{U}}_i^*$.

Finally, the adaptive LASSO coefficient estimates at location \mathbf{s}_i are found by using the `lars` algorithm (Efron et al., 2004) to model \mathbf{V}_i^* as a function of $\check{\mathbf{U}}_i^*$.

3.2. Tuning parameter selection

At each location \mathbf{s}_i , it is necessary to select the LASSO tuning parameter λ_i . To compare different values of λ_i , we propose a locally-weighted version of the Akaike information criterion (AIC (Akaike, 1974)) which we call the local AIC, or AIC_{loc} . The local AIC is calculated by adding a penalty to the local likelihood, with the sum of the weights around \mathbf{s}_i , $\sum_{i'=1}^n w_{ii'}$, playing the role of the sample size and the “degrees of freedom” (df_i) at \mathbf{s}_i given by the number of nonzero coefficients in $\boldsymbol{\beta}_i$ (Zou et al., 2007).

$$\text{AIC}_{\text{loc},i} = -2 \sum_{i'=1}^n \ell_{ii'} + 2\text{df}_i \quad (13)$$

$$= -2 \times \sum_{i'=1}^n \log \left\{ (2\pi\hat{\sigma}_i^2)^{-1/2} \exp \left[-\frac{1}{2}\hat{\sigma}_i^{-2} (y_{i'} - \mathbf{x}_{i'}'\hat{\boldsymbol{\beta}}_i)^2 \right] \right\}^{w_{ii'}} + 2\text{df}_i \quad (14)$$

$$= \sum_{i'=1}^n w_{ii'} \left\{ \log(2\pi) + \log \hat{\sigma}_i^2 + \hat{\sigma}_i^{-2} (y_{i'} - \mathbf{x}_{i'}'\hat{\boldsymbol{\beta}}_i)^2 \right\} + 2\text{df}_i \quad (15)$$

$$= \hat{\sigma}_i^{-2} \sum_{i'=1}^n w_{ii'} (y_{i'} - \mathbf{x}_{i'}'\hat{\boldsymbol{\beta}}_i)^2 + 2\text{df}_i + C_i \quad (16)$$

Where the estimated local variance $\hat{\sigma}_i^2$ is the variance estimate from the unpenalized local model (Zou et al., 2007), so C_i does not depend on the choice of tuning parameter and can be ignored.

Wheeler (2009) proposed selecting the tuning parameter for the LASSO at location \mathbf{s}_i to minimize the jackknife prediction error $|y_i - \hat{y}_i^{(i)}|$. Because the jackknife prediction error is undefined everywhere except for at observation locations, this choice restricts coefficient estimation to occur at the locations where data has been observed. By contrast, the local AIC can be calculated at any location where we can calculate the local likelihood. As a practical matter this allows for variable selection and coefficient surface estimation to be done at locations where no data was observed (interpolation) and for imputation of missing values of the response variable.

3.3. Bandwidth selection

The bandwidth parameter is global and so we need a global statistic for comparing prospective bandwidths. The objective minimized by GWL is:

$$\sum_{i'=1}^n w_{ii'} (y_{i'} - \mathbf{x}'_{i'} \boldsymbol{\beta}_i)^2 + \sum_{j=1}^p \lambda_{ij} \beta_{ij} \quad (17)$$

Where $\lambda_{ij}, j = 1, \dots, p$ are penalties from the adaptive LASSO (Zou, 2006). Taking the derivatives with respect to $\boldsymbol{\beta}$ and setting to zero, we see that

$$\hat{\boldsymbol{\beta}}_{i,\text{GWL}} = (\mathbf{X}' \mathbf{W}_i \mathbf{X})^{-1} \mathbf{X}' \mathbf{W}_i \mathbf{Y} - \frac{1}{2} (\mathbf{X}' \mathbf{W}_i \mathbf{X})^{-1} \boldsymbol{\lambda}_i \quad (18)$$

$$\hat{y}_i = \mathbf{x}'_i \hat{\boldsymbol{\beta}}_{i,\text{GWL}} = \mathbf{x}'_i (\mathbf{X}' \mathbf{W}_i \mathbf{X})^{-1} \mathbf{X}' \mathbf{W}_i \mathbf{Y} - \frac{1}{2} \mathbf{x}'_i (\mathbf{X}' \mathbf{W}_i \mathbf{X})^{-1} \boldsymbol{\lambda}_i \quad (19)$$

Unlike in the case of ordinary geographically-weighted regression, the fitted values $\hat{\mathbf{Y}}$ are not a linear combination of the observations \mathbf{Y} . Because GWL is not a linear smoother it is not possible to calculate the AIC as in Fotheringham et al. (2002) (Zou, 2006). We propose a statistic called the total AIC (AIC_{tot}) for the purpose of selecting the bandwidth parameter. Because of the kernel weights and the application of the adaptive LASSO, the sample size and the degrees of freedom are different at each location. The total AIC is found by taking the sum over all of the observed data:

$$\text{AIC}_{\text{tot}} = -2 \times \sum_{i=1}^n \ell_i + 2 \times \text{df} \quad (20)$$

$$= \sum_{i=1}^n \left\{ \log \hat{\sigma}_i^2 + \hat{\sigma}_i^{-2} \left(y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_i \right)^2 \right\} + 2 \times \text{df} \quad (21)$$

What remains is to calculate df, the number of degrees of freedom used by the model. Classical GWR, as developed in Loader (1999) and Fotheringham et al. (2002) calculates df using the trace of the “hat” matrix, but because the GWL is not a linear smoother, there is no “hat” matrix associated with GWL. Instead, notice that df can be pulled into the summation in (21):

$$\text{df} = \sum_{i=1}^n (n^{-1} \text{df}) \quad (22)$$

Now, because we are considering the sum of local weights to be the sample size for the local models, we estimate df by $\sum_{i=1}^n \left\{ \left(\sum_{i'=1}^n w_{ii'} \right)^{-1} \text{df}_i \right\}$, and the total AIC is then:

$$\text{AIC}_{\text{tot}} = \sum_{i=1}^n \left\{ \log \hat{\sigma}_i^2 + \hat{\sigma}_i^{-2} \left(y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_i \right)^2 + 2 \times \left(\sum_{i'=1}^n w_{ii'} \right)^{-1} \text{df}_i \right\} \quad (23)$$

The bandwidth that minimizes (23) is found by a line search.

3.4. Confidence interval construction

Confidence intervals for the GWL’s coefficient estimates can be calculated either by the bootstrap (Efron and Tibshirani, 1986) or by exploiting an assumption of normally-distributed residuals. Then the, e.g., 95% confidence interval for each regression coefficient is defined by the (2.5, 97.5)

percentiles of the coefficient estimates from the bootstrap replicates.

To compute coefficient confidence intervals via the bootstrap, the observations with non-zero geographic weights are resampled uniformly with replacement for each of n_B bootstrap replicates. For each bootstrap replicate, the GWL is used to estimate regression coefficients. The local likelihood of the bootstrap replicates may be different from that of the original sample, so the adaptive LASSO tuning parameter may differ for each bootstrap replicate. Since the GWL is applied independently to each bootstrap replicate, the variables selected by GWL may be different for each replicate.

Unshrunk coefficient estimates are found by using the GWL at each location for variable selection only and then estimating the coefficients for the selected variables by weighted least squares. An unshrunk bootstrap confidence interval is found by estimating the unshrunk coefficients for each of the n_B bootstrap replicates and then calculating the percentiles as above.

A third way to estimate the coefficient confidence intervals is to use the GWL for variable selection only and then to use weighted least squares for both coefficient estimation and confidence interval construction:

$$\hat{\text{se}}_{\beta_i} = \left(\tilde{\mathbf{X}}_i' \mathbf{W}_i \tilde{\mathbf{X}}_i \right)^{-1} \tilde{\mathbf{X}}_i' \mathbf{W}_i \mathbf{Y} \quad (24)$$

where $\tilde{\mathbf{X}}_i$ is the model matrix including only those variables that are selected by GWL at location i .

4. Simulation

4.1. Simulation setup

A simulation study was conducted to assess some finite-sample properties of the method described in Sections 2-3. Data was simulated on $[0, 1] \times [0, 1]$, which was divided into a 30×30 grid. Each of $p = 5$ covariates Z_1, \dots, Z_p was simulated by a Gaussian random field (GRF) with mean zero and exponential spatial covariance $Cov(Z_{ji}, Z_{ji'}) = \sigma_z^2 \exp(-\tau_z^{-1} \delta_{ii'})$ where $\sigma_z^2 = 1$ is the variance, τ_z is the range parameter, and $\delta_{ii'}$ is the Euclidean distance $\|\mathbf{s}_i - \mathbf{s}_{i'}\|_2$. Correlation was induced between the covariates by multiplying the \mathbf{Z} matrix by \mathbf{R} , where \mathbf{R} is the Cholesky decomposition of the covariance matrix $\Sigma = \mathbf{R}'\mathbf{R}$. The covariance matrix Σ is a 5×5 matrix that has ones on the diagonal and ρ for all off-diagonal entries, where ρ is the between-covariate correlation.

The simulated response is $y_i = \mathbf{z}_i' \boldsymbol{\beta}_i + \epsilon_i$ for $i = 1, \dots, 900$ where the vector of additive errors $\boldsymbol{\epsilon}$ is generated from a GRF with spatial covariance $Cov(\epsilon_i, \epsilon_{i'}) = \sigma_\epsilon^2 \exp(-\tau_\epsilon^{-1} \delta_{ii'})$ where $\sigma_\epsilon^2 = 1$.

The simulated data include the output y and five covariates Z_1, \dots, Z_5 . The true data-generating model uses only Z_1 , so Z_2, \dots, Z_5 are included to test the variable-selection properties of GWL. The coefficient surface of β_1 is described by the “step” function (Figure 1):

$$\beta_1(s) = \begin{cases} 0 & \text{if } s_y < 0.4 \\ 5(s_y - 0.4) & \text{if } 0.4 \leq s_y < 0.6 \\ 1 & \text{o.w.} \end{cases} \quad (25)$$

In order to evaluate the performance of GWL under a range of conditions, the data was simulated under 18 different settings (Table 1): high (0.1) and low (0.03) levels of τ_z , the autoregression range parameter for the covariate GRFs; three levels (0, 0.5, 0.8) of between-covariate correlation ρ ; and three levels (0, 0.03, 0.1) of the autoregression range parameter τ_ϵ for the error-term GRF $\boldsymbol{\epsilon}$. Each

case was simulated 100 times.

The performance of GWL was compared to oracular GWR (O-GWR), which is ordinary GWR with “oracular” variable selection, meaning that exactly the correct set of predictors was used to fit the GWR model at each location in the simulation.

[Table 1 about here.]

[Figure 1 about here.]

4.2. Results

Results from the simulation were summarized at five locations on the simulated grid. The five key locations were chosen because they represent interesting regions of the β_1 coefficient surface. The results of variable selection (Tables 2 - 6) and coefficient estimation (Tables 7 - 11) are presented in the tables below.

4.3. Tables

4.3.1. Selection

[Table 2 about here.]

[Table 3 about here.]

[Table 4 about here.]

[Table 5 about here.]

[Table 6 about here.]

4.3.2. Estimation

[Table 7 about here.]

[Table 8 about here.]

[Table 9 about here.]

[Table 10 about here.]

[Table 11 about here.]

5. Data analysis

5.1. Census poverty data

We present the following analysis to demonstrate one possible application of the geographically-weighted LASSO in a linear regression context. We use county-level data from the U.S. Census Bureau to select the social and demographic variables that are important predictors of the county-level poverty rate in the upper midwest, and to estimate the coefficients associated with these predictors. Data are from six censuses - the decennial censuses from 1960 to 2000, and from the American Community Survey in 2006. Selection and estimation are done for each census individually (no attempt is made here to borrow strength across years). The outcome of interest (poverty rate) is a proportion and so takes values on $[0, 1]$, but to demonstrate the GWL in a linear regression context, we model the logit-transformed poverty rate. Our data set covers all counties in the states of Minnesota, Iowa, Wisconsin, Illinois, Indiana, and Michigan. The potential predictors are described in Table 12.

[Table 12 about here.]

5.2. Figures

The coefficient estimates are plotted on maps of the upper midwest in Figures 2 - 7. It is immediately apparent that the estimated coefficient surfaces are non-constant for most variables.

[Figure 2 about here.]

[Figure 3 about here.]

[Figure 4 about here.]

[Figure 5 about here.]

[Figure 6 about here.]

[Figure 7 about here.]

6. References

References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control* 19(6), 716–723.
- Antoniadas, A., I. Gijbels, and A. Verhasselt (2012). Variable selection in varying-coefficient models using p-splines. *Journal of Computational and Graphical Statistics* 21(3), 638–661.
- Breiman, L. (1995). Better subset regression using the nonnegative garrote. *Technometrics* 51, 373–384.
- Efron, B., T. Hastie, I. Johnstone, and R. Tibshirani (2004). Least angle regression. *Annals of Statistics* 32(2), 407–499.

- Efron, B. and R. Tibshirani (1986). Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. *Statistical Science* 1(1), 54–75.
- Fan, J. and W. Zhang (1999). Statistical estimation in varying coefficient models. *The Annals of Statistics* 27(5), 1491–1518.
- Fotheringham, A., C. Brunsdon, and M. Charlton (2002). *Geographically weighted regression: the analysis of spatially varying relationships*. Wiley.
- Hastie, T. and C. Loader (1993). Local regression: automatic kernel carpentry. *Statistical Science* 8(2), 120–143.
- Hastie, T. and R. Tibshirani (1993). Varying-coefficient models. *Journal of the Royal Statistical Society. Series B (Methodological)* 55(4), pp. 757–796.
- Loader, C. (1999). *Local regression and likelihood*. Springer New York.
- Wang, L., H. Li, and J. Z. Huang (2008). Variable selection in nonparametric varying-coefficient models for analysis of repeated measurements. *Journal of the American Statistical Association* 103(484), 1556–1569.
- Wang, N., C.-L. Mei, and X.-D. Yan (2008). Local linear estimation of spatially varying coefficient models: an improvement on the geographically weighted regression technique. *Environment and Planning A* 40, 986–1005.
- Wheeler, D. C. (2009). Simultaneous coefficient penalization and model selection in geographically weighted regression: the geographically weighted lasso. *Environment and Planning A* 41, 722–742.
- Wood, S. (2006). *Generalized additive models: an introduction with R*. Texts in statistical science. Chapman & Hall/CRC.

Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* 101(476), 1418–1429.

Zou, H., T. Hastie, and R. Tibshirani (2007). On the “degrees of freedom” of the lasso. *Annals of Statistics* 35(5), 2173–2192.

List of Figures

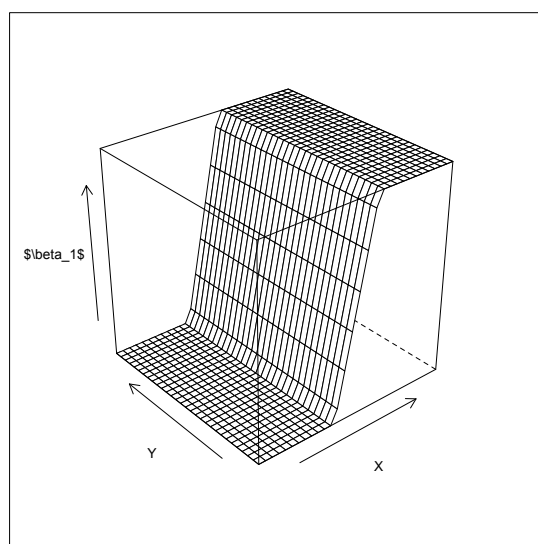


Figure 1: The actual β_1 coefficient surface used in the simulation.

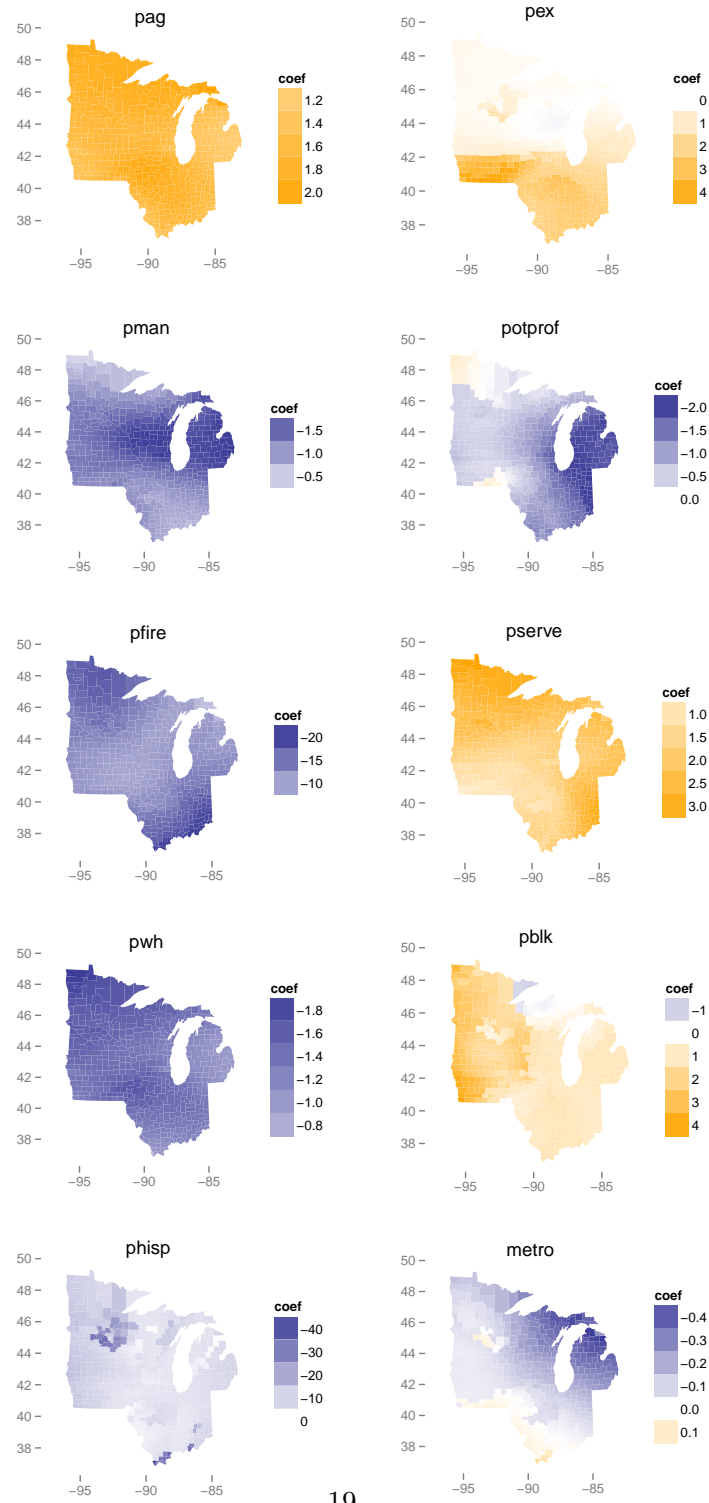


Figure 2: Estimated coefficient surfaces for the 1960 census.

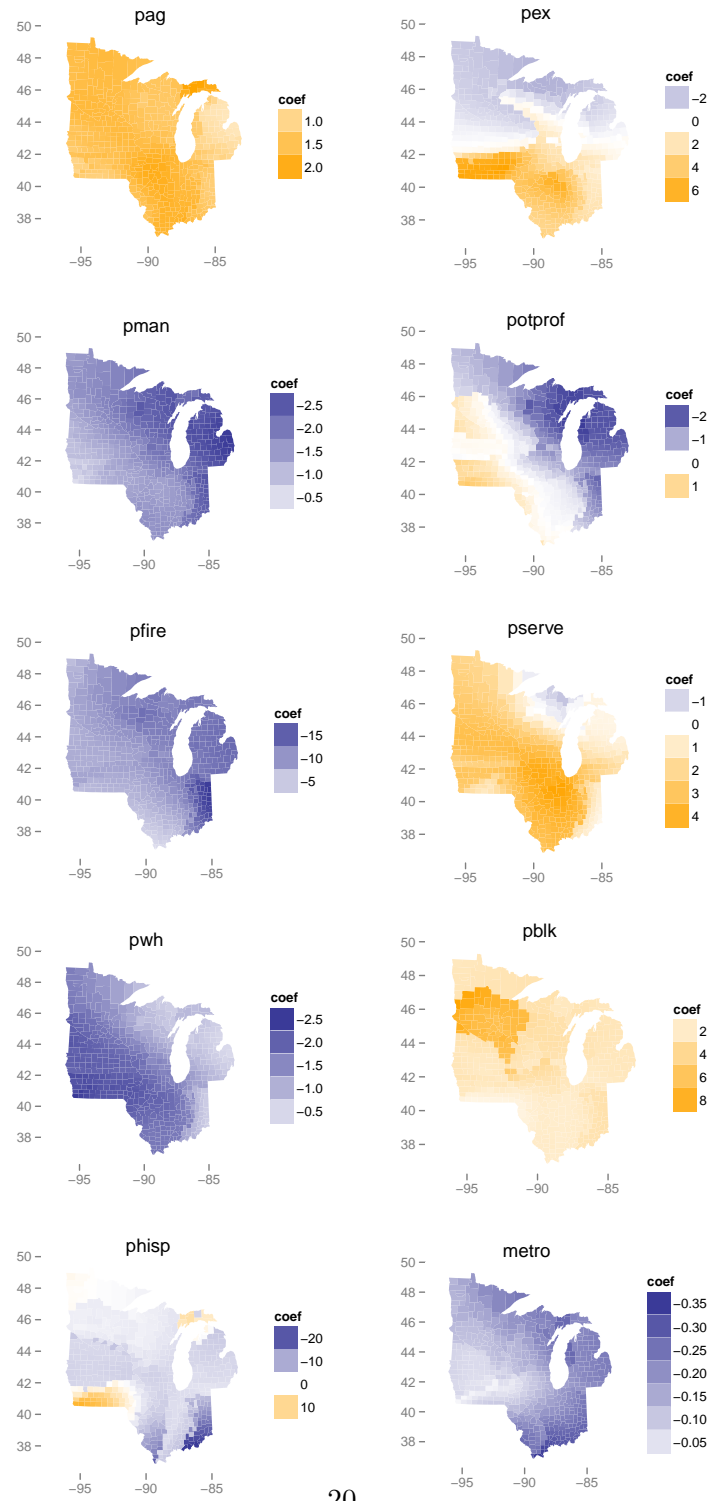


Figure 3: Estimated coefficient surfaces for the 1970 census.

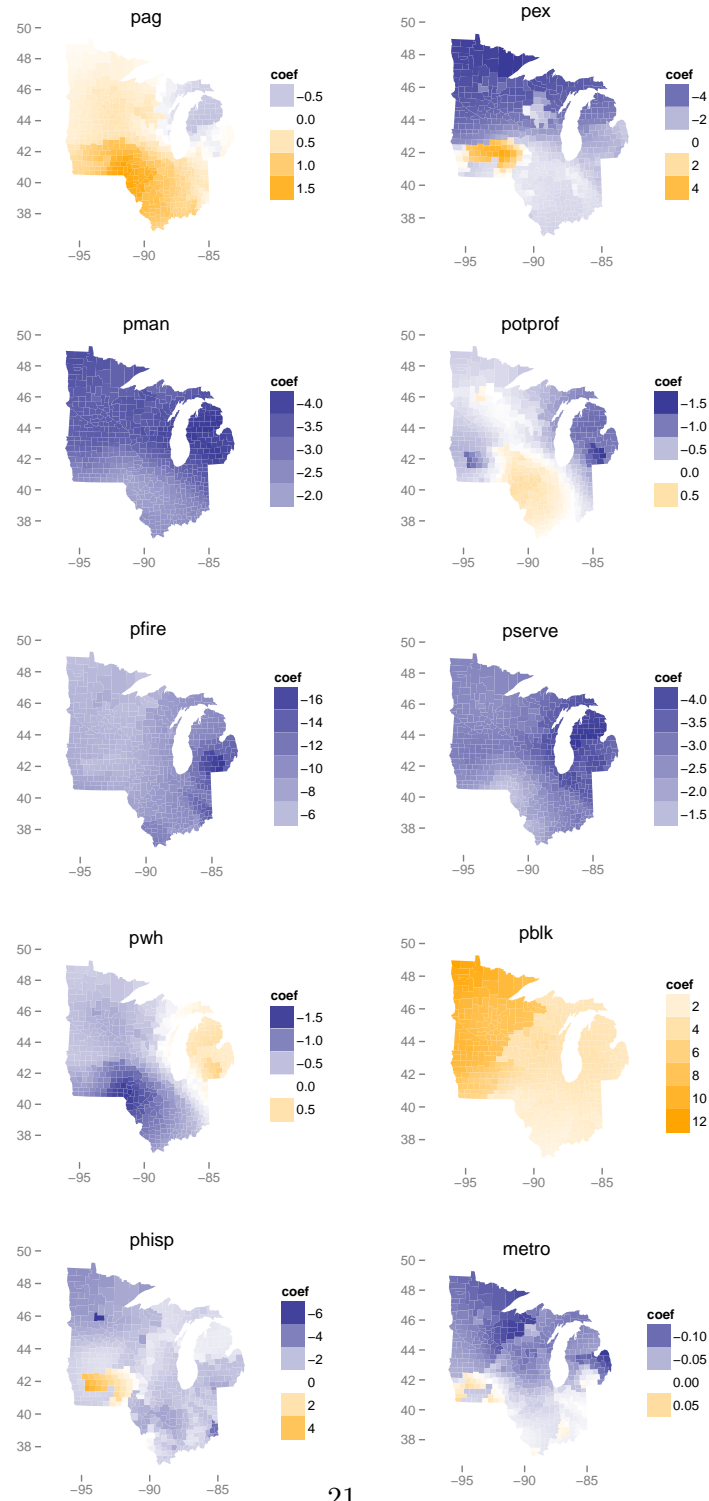


Figure 4: Estimated coefficient surfaces for the 1980 census.

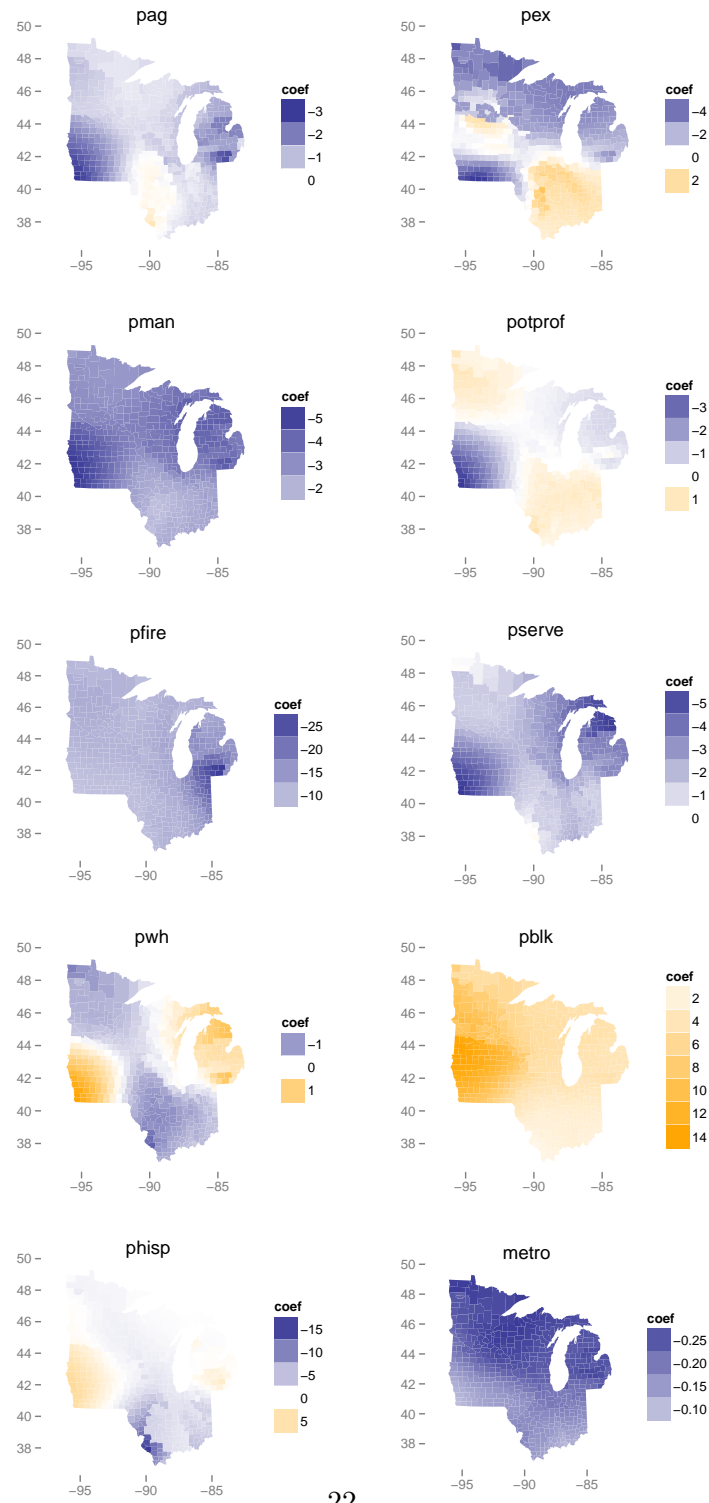


Figure 5: Estimated coefficient surfaces for the 1990 census.

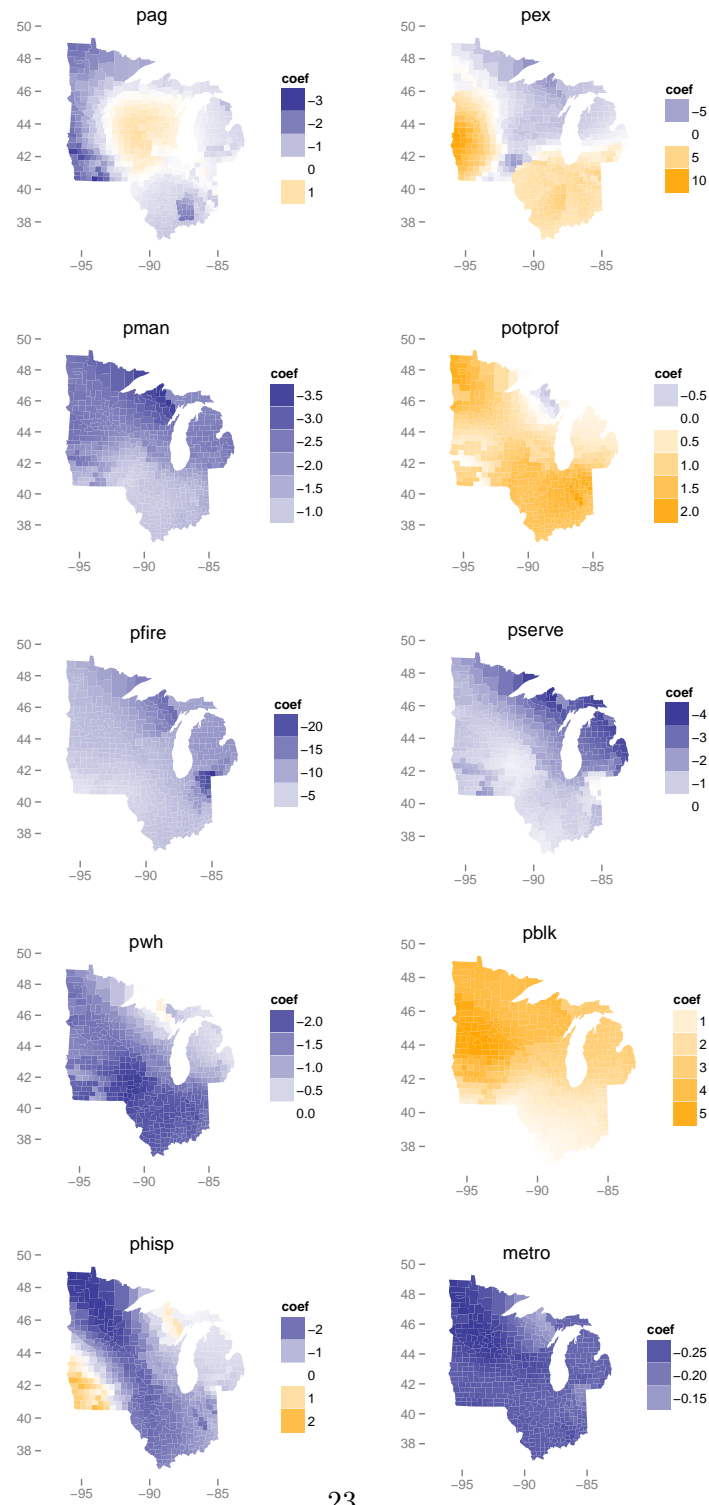


Figure 6: Estimated coefficient surfaces for the 2000 census.

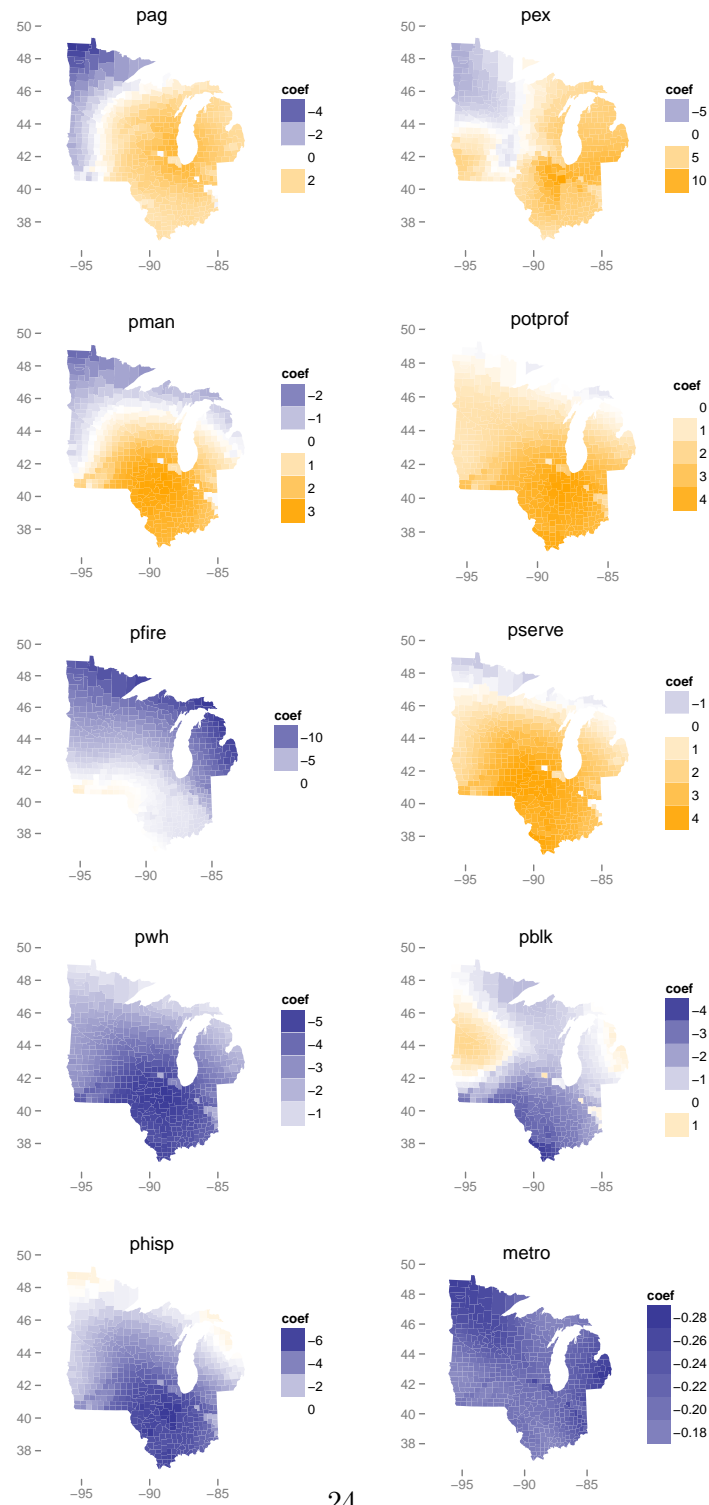


Figure 7: Estimated coefficient surfaces for the 2006 census.

List of Tables

Setting	τ_z	ρ	τ_ϵ
1	0.03	0.00	0.00
2	0.03	0.00	0.03
3	0.03	0.00	0.10
4	0.03	0.50	0.00
5	0.03	0.50	0.03
6	0.03	0.50	0.10
7	0.03	0.80	0.00
8	0.03	0.80	0.03
9	0.03	0.80	0.10
10	0.10	0.00	0.00
11	0.10	0.00	0.03
12	0.10	0.00	0.10
13	0.10	0.50	0.00
14	0.10	0.50	0.03
15	0.10	0.50	0.10
16	0.10	0.80	0.00
17	0.10	0.80	0.03
18	0.10	0.80	0.10

Table 1: Simulation parameters for each setting.

ρ	τ_x	τ_σ	β_1	$\beta_2 - \beta_5$
0	0.03	0	1.00	0.21
		0.03	0.99	0.34
		0.1	1.00	0.44
	0.1	0	1.00	0.20
		0.03	0.99	0.38
		0.1	0.97	0.59
0.5	0.03	0	1.00	0.26
		0.03	1.00	0.29
		0.1	1.00	0.47
	0.1	0	1.00	0.26
		0.03	0.98	0.39
		0.1	0.98	0.60
0.8	0.03	0	0.99	0.30
		0.03	0.94	0.40
		0.1	0.98	0.45
	0.1	0	1.00	0.34
		0.03	0.93	0.46
		0.1	0.94	0.61

Table 2: Selection frequency for the indicated variables at location 1

ρ	τ_x	τ_σ	β_1	$\beta_2 - \beta_5$
0	0.03	0	1.00	0.14
		0.03	1.00	0.24
		0.1	1.00	0.38
	0.1	0	1.00	0.29
		0.03	1.00	0.40
		0.1	1.00	0.71
0.5	0.03	0	1.00	0.17
		0.03	1.00	0.29
		0.1	1.00	0.40
	0.1	0	1.00	0.27
		0.03	1.00	0.42
		0.1	1.00	0.70
0.8	0.03	0	1.00	0.19
		0.03	0.96	0.33
		0.1	1.00	0.40
	0.1	0	1.00	0.37
		0.03	0.93	0.42
		0.1	0.96	0.63

Table 3: Selection frequency for the indicated variables at location 2

ρ	τ_x	τ_σ	β_1	$\beta_2 - \beta_5$
0	0.03	0	1.00	0.17
		0.03	1.00	0.35
		0.1	1.00	0.50
	0.1	0	1.00	0.20
		0.03	0.97	0.40
		0.1	1.00	0.62
0.5	0.03	0	1.00	0.21
		0.03	0.99	0.37
		0.1	1.00	0.43
	0.1	0	1.00	0.23
		0.03	0.97	0.38
		0.1	1.00	0.71
0.8	0.03	0	0.99	0.27
		0.03	0.96	0.32
		0.1	1.00	0.48
	0.1	0	0.99	0.29
		0.03	0.96	0.42
		0.1	0.97	0.70

Table 4: Selection frequency for the indicated variables at location 3

ρ	τ_x	τ_σ	β_1	$\beta_2 - \beta_5$
0	0.03	0	0.77	0.15
		0.03	0.54	0.22
		0.1	0.85	0.33
	0.1	0	0.90	0.27
		0.03	0.74	0.39
		0.1	0.80	0.60
0.5	0.03	0	0.71	0.18
		0.03	0.54	0.25
		0.1	0.80	0.38
	0.1	0	0.85	0.19
		0.03	0.60	0.39
		0.1	0.67	0.56
0.8	0.03	0	0.70	0.21
		0.03	0.35	0.17
		0.1	0.72	0.35
	0.1	0	0.85	0.26
		0.03	0.45	0.33
		0.1	0.68	0.61

Table 5: Selection frequency for the indicated variables at location 4

ρ	τ_x	τ_σ	β_1	$\beta_2 - \beta_5$
0	0.03	0	0.05	0.09
		0.03	0.19	0.18
		0.1	0.35	0.34
	0.1	0	0.05	0.07
		0.03	0.26	0.21
		0.1	0.55	0.60
0.5	0.03	0	0.09	0.10
		0.03	0.16	0.22
		0.1	0.30	0.30
	0.1	0	0.09	0.09
		0.03	0.18	0.20
		0.1	0.45	0.51
0.8	0.03	0	0.08	0.05
		0.03	0.13	0.21
		0.1	0.36	0.33
	0.1	0	0.05	0.04
		0.03	0.21	0.25
		0.1	0.47	0.55

Table 6: Selection frequency for the indicated variables at location 5

ρ	τ_x	τ_σ	GWL			GWL-U			Oracle		
			MSE	bias	var	MSE	bias	var	MSE	bias	var
0		0	0.048	0.122	0.034	0.030	<i>0.002</i>	0.030	<i>0.032</i>	-0.001	<i>0.032</i>
	0.03	0.03	<i>0.071</i>	0.090	0.063	0.064	-0.015	<i>0.064</i>	0.138	<i>-0.039</i>	0.138
		0.1	<i>0.044</i>	0.084	<i>0.038</i>	0.028	0.006	0.028	0.258	<i>0.032</i>	0.259
		0	0.049	0.119	<i>0.035</i>	0.027	-0.001	0.027	<i>0.041</i>	<i>-0.009</i>	0.041
	0.1	0.03	<i>0.088</i>	0.120	0.074	0.083	<i>-0.004</i>	<i>0.084</i>	0.221	0.000	0.223
		0.1	<i>0.109</i>	<i>0.081</i>	<i>0.104</i>	0.093	0.027	0.093	0.508	0.089	0.505
0.5		0	0.067	0.158	0.042	<i>0.042</i>	<i>0.055</i>	<i>0.039</i>	0.033	-0.008	0.033
	0.03	0.03	<i>0.105</i>	0.151	0.083	0.085	0.029	<i>0.085</i>	0.110	<i>-0.036</i>	0.110
		0.1	<i>0.062</i>	0.098	<i>0.053</i>	0.045	<i>0.028</i>	0.045	0.285	0.015	0.288
		0	<i>0.059</i>	0.147	0.038	0.046	<i>0.052</i>	<i>0.044</i>	0.074	0.029	0.074
	0.1	0.03	<i>0.134</i>	0.161	0.109	0.123	<i>0.058</i>	<i>0.121</i>	0.208	-0.043	0.209
		0.1	<i>0.140</i>	0.061	<i>0.138</i>	0.114	0.005	0.115	0.620	<i>-0.044</i>	0.624
0.8		0	0.128	0.220	<i>0.080</i>	<i>0.116</i>	<i>0.147</i>	0.095	0.031	0.015	0.031
	0.03	0.03	<i>0.255</i>	0.198	<i>0.218</i>	0.270	<i>0.096</i>	0.263	0.097	-0.026	0.097
		0.1	<i>0.131</i>	0.122	0.117	0.127	0.049	<i>0.125</i>	0.273	<i>-0.100</i>	0.266
		0	<i>0.117</i>	0.208	<i>0.074</i>	0.119	<i>0.132</i>	0.103	0.042	0.016	0.042
	0.1	0.03	<i>0.424</i>	0.133	<i>0.411</i>	0.521	<i>0.044</i>	0.524	0.212	-0.013	0.214
		0.1	0.290	<i>0.033</i>	0.291	<i>0.300</i>	0.009	<i>0.303</i>	0.680	0.083	0.680

Table 7: MSE, bias, and variance of estimates for β_1 at location 1 (**minimum**, *next best*).

ρ	τ_x	τ_σ	GWL			GWL-U			Oracle		
			MSE	bias	var	MSE	bias	var	MSE	bias	var
0		0	0.083	0.265	0.013	<i>0.048</i>	<i>0.190</i>	<i>0.013</i>	0.042	0.171	0.013
	0.03	0.03	0.081	0.244	<i>0.022</i>	0.042	<i>0.146</i>	0.021	<i>0.052</i>	0.052	0.049
		0.1	<i>0.089</i>	0.280	<i>0.011</i>	0.062	<i>0.227</i>	0.011	0.097	0.048	0.096
		0	0.087	0.277	<i>0.011</i>	<i>0.060</i>	<i>0.223</i>	0.011	0.053	0.192	0.016
	0.1	0.03	<i>0.103</i>	0.262	0.035	0.071	<i>0.188</i>	<i>0.036</i>	0.105	0.098	0.097
		0.1	<i>0.095</i>	0.234	<i>0.041</i>	0.071	<i>0.197</i>	0.032	0.134	0.013	0.135
0.5		0	0.102	0.292	0.016	<i>0.065</i>	<i>0.221</i>	<i>0.016</i>	0.045	0.177	0.014
	0.03	0.03	0.093	0.252	0.029	<i>0.064</i>	<i>0.169</i>	<i>0.036</i>	0.043	0.072	0.038
		0.1	0.104	0.280	<i>0.025</i>	0.075	<i>0.227</i>	0.024	<i>0.076</i>	-0.002	0.077
		0	0.110	0.307	0.016	<i>0.082</i>	<i>0.255</i>	0.017	0.059	0.206	<i>0.016</i>
	0.1	0.03	0.115	0.262	0.048	0.092	<i>0.204</i>	<i>0.051</i>	<i>0.104</i>	0.090	0.097
		0.1	<i>0.093</i>	0.187	0.059	0.087	<i>0.162</i>	<i>0.061</i>	0.218	0.057	0.217
0.8		0	0.126	0.322	<i>0.023</i>	<i>0.094</i>	<i>0.255</i>	0.030	0.041	0.172	0.011
	0.03	0.03	0.166	0.317	<i>0.066</i>	<i>0.137</i>	<i>0.255</i>	0.073	0.047	0.106	0.036
		0.1	0.117	0.286	0.035	<i>0.098</i>	<i>0.240</i>	<i>0.041</i>	0.074	0.026	0.074
		0	0.132	0.319	<i>0.031</i>	<i>0.120</i>	<i>0.283</i>	0.040	0.049	0.172	0.019
	0.1	0.03	<i>0.212</i>	0.349	<i>0.091</i>	0.216	<i>0.303</i>	0.125	0.097	0.114	0.085
		0.1	0.207	0.272	0.135	<i>0.203</i>	<i>0.228</i>	<i>0.153</i>	0.168	0.072	0.164

Table 8: MSE, bias, and variance of estimates for β_1 at location 2 (**minimum**, *next best*).

ρ	τ_x	τ_σ	GWL			GWL-U			Oracle		
			MSE	bias	var	MSE	bias	var	MSE	bias	var
0		0	0.028	0.102	0.018	0.014	<i>0.003</i>	0.014	<i>0.015</i>	-0.001	<i>0.015</i>
	0.03	0.03	<i>0.042</i>	0.095	<i>0.033</i>	0.030	0.001	0.030	0.093	<i>0.023</i>	0.094
		0.1	<i>0.034</i>	0.071	<i>0.029</i>	0.019	<i>0.013</i>	0.019	0.168	-0.004	0.169
		0	<i>0.022</i>	0.075	<i>0.017</i>	0.016	-0.012	0.016	0.023	<i>-0.018</i>	0.023
	0.1	0.03	<i>0.083</i>	0.113	0.071	0.078	0.022	<i>0.078</i>	0.157	<i>-0.037</i>	0.157
		0.1	<i>0.078</i>	0.080	<i>0.073</i>	0.049	0.004	0.049	0.250	<i>0.064</i>	0.248
0.5		0	0.036	0.125	<i>0.021</i>	<i>0.028</i>	<i>0.044</i>	0.026	0.018	0.002	0.018
	0.03	0.03	<i>0.071</i>	0.149	0.049	0.053	<i>0.060</i>	<i>0.050</i>	0.074	0.034	0.073
		0.1	<i>0.036</i>	0.070	<i>0.031</i>	0.026	0.017	0.026	0.173	<i>0.023</i>	0.174
		0	0.023	0.101	0.013	0.017	<i>0.033</i>	<i>0.016</i>	<i>0.020</i>	-0.005	0.020
	0.1	0.03	<i>0.108</i>	0.133	0.091	0.104	0.045	<i>0.103</i>	0.148	<i>0.064</i>	0.146
		0.1	<i>0.101</i>	0.010	<i>0.102</i>	0.094	-0.023	0.094	0.469	<i>0.017</i>	0.474
0.8		0	0.073	0.182	<i>0.041</i>	<i>0.062</i>	<i>0.132</i>	0.045	0.015	0.032	0.014
	0.03	0.03	0.162	0.206	<i>0.120</i>	<i>0.152</i>	<i>0.116</i>	0.140	0.075	-0.010	0.076
		0.1	<i>0.091</i>	0.086	<i>0.084</i>	0.079	<i>0.036</i>	0.078	0.146	-0.028	0.147
		0	<i>0.057</i>	0.132	<i>0.040</i>	0.061	<i>0.086</i>	0.054	0.019	-0.023	0.019
	0.1	0.03	<i>0.182</i>	0.205	<i>0.141</i>	0.185	<i>0.102</i>	0.177	0.128	-0.027	0.129
		0.1	<i>0.216</i>	0.050	<i>0.216</i>	0.210	0.012	0.212	0.326	<i>0.029</i>	0.328

Table 9: MSE, bias, and variance of estimates for β_1 at location 3 (**minimum**, *next best*).

ρ	τ_x	τ_σ	GWL			GWL-U			Oracle		
			MSE	bias	var	MSE	bias	var	MSE	bias	var
0		0	0.024	-0.101	<i>0.014</i>	0.044	<i>-0.155</i>	0.020	<i>0.037</i>	-0.159	0.012
	0.03	0.03	0.023	<i>-0.060</i>	0.020	<i>0.036</i>	-0.094	<i>0.028</i>	0.046	-0.038	0.045
		0.1	0.032	<i>-0.136</i>	0.014	<i>0.047</i>	-0.176	<i>0.017</i>	0.091	-0.005	0.092
		0	0.042	-0.162	<i>0.015</i>	0.065	-0.220	0.017	<i>0.049</i>	<i>-0.188</i>	0.013
	0.1	0.03	0.031	<i>-0.114</i>	0.019	<i>0.057</i>	-0.161	<i>0.031</i>	0.100	-0.061	0.097
		0.1	0.055	<i>-0.157</i>	0.031	<i>0.064</i>	-0.180	<i>0.031</i>	0.210	-0.042	0.210
0.5		0	0.023	-0.087	<i>0.015</i>	0.040	<i>-0.130</i>	0.023	<i>0.038</i>	-0.161	0.012
	0.03	0.03	0.024	-0.061	0.021	<i>0.039</i>	-0.093	<i>0.031</i>	0.041	<i>-0.085</i>	0.035
		0.1	0.037	<i>-0.128</i>	0.021	<i>0.052</i>	-0.165	<i>0.025</i>	0.059	-0.023	0.059
		0	0.041	-0.151	<i>0.018</i>	0.067	-0.208	0.024	<i>0.055</i>	<i>-0.198</i>	0.016
	0.1	0.03	0.046	<i>-0.095</i>	0.037	<i>0.063</i>	-0.124	<i>0.048</i>	0.090	-0.058	0.088
		0.1	0.078	<i>-0.137</i>	0.060	<i>0.089</i>	-0.149	<i>0.068</i>	0.153	0.038	0.153
0.8		0	0.024	-0.089	<i>0.016</i>	<i>0.036</i>	<i>-0.118</i>	0.022	0.041	-0.178	0.009
	0.03	0.03	0.025	-0.014	0.025	0.047	-0.041	0.045	<i>0.038</i>	<i>-0.037</i>	<i>0.038</i>
		0.1	0.031	<i>-0.101</i>	0.021	<i>0.050</i>	-0.142	<i>0.030</i>	0.079	-0.040	0.078
		0	0.044	-0.148	<i>0.023</i>	0.062	<i>-0.186</i>	0.027	<i>0.055</i>	-0.200	0.015
	0.1	0.03	0.056	-0.042	0.055	0.101	<i>-0.066</i>	0.097	<i>0.081</i>	-0.069	<i>0.077</i>
		0.1	<i>0.157</i>	<i>-0.188</i>	0.122	0.168	-0.215	<i>0.123</i>	0.155	-0.028	0.156

Table 10: MSE, bias, and variance of estimates for β_1 at location 4 (**minimum**, *next best*).

ρ	τ_x	τ_σ	GWL			GWL-U			Oracle		
			MSE	bias	var	MSE	bias	var	MSE	bias	var
0		0	<i>0.002</i>	<i>0.003</i>	<i>0.002</i>	0.004	0.004	0.004	0.000	0.000	0.000
	0.03	0.03	<i>0.029</i>	0.024	<i>0.029</i>	0.035	<i>0.024</i>	0.035	0.000	0.000	0.000
		0.1	<i>0.017</i>	0.008	<i>0.017</i>	0.018	<i>0.003</i>	0.018	0.000	0.000	0.000
		0	<i>0.002</i>	<i>-0.003</i>	<i>0.002</i>	0.005	-0.004	0.005	0.000	0.000	0.000
	0.1	0.03	<i>0.039</i>	-0.007	<i>0.040</i>	0.054	<i>0.000</i>	0.054	0.000	0.000	0.000
		0.1	0.079	0.061	0.076	<i>0.069</i>	<i>0.059</i>	<i>0.066</i>	0.000	0.000	0.000
0.5		0	<i>0.010</i>	-0.010	<i>0.010</i>	0.016	<i>-0.009</i>	0.016	0.000	0.000	0.000
	0.03	0.03	<i>0.022</i>	<i>0.021</i>	<i>0.022</i>	0.035	0.024	0.035	0.000	0.000	0.000
		0.1	0.025	<i>0.017</i>	0.025	<i>0.024</i>	0.026	<i>0.023</i>	0.000	0.000	0.000
		0	<i>0.008</i>	<i>-0.018</i>	<i>0.007</i>	0.018	-0.029	0.018	0.000	0.000	0.000
	0.1	0.03	<i>0.030</i>	<i>0.000</i>	<i>0.030</i>	0.040	0.011	0.040	0.000	0.000	0.000
		0.1	0.057	-0.036	0.057	<i>0.057</i>	<i>-0.034</i>	<i>0.056</i>	0.000	0.000	0.000
0.8		0	<i>0.007</i>	<i>-0.004</i>	<i>0.008</i>	0.025	-0.008	0.026	0.000	0.000	0.000
	0.03	0.03	<i>0.041</i>	<i>-0.024</i>	<i>0.040</i>	0.053	-0.038	0.052	0.000	0.000	0.000
		0.1	<i>0.059</i>	<i>-0.001</i>	<i>0.059</i>	0.065	-0.008	0.066	0.000	0.000	0.000
		0	<i>0.004</i>	0.003	<i>0.004</i>	0.008	<i>0.002</i>	0.008	0.000	0.000	0.000
	0.1	0.03	<i>0.116</i>	<i>-0.035</i>	<i>0.116</i>	0.173	-0.056	0.172	0.000	0.000	0.000
		0.1	0.160	0.069	0.157	<i>0.132</i>	<i>0.028</i>	<i>0.133</i>	0.000	0.000	0.000

Table 11: MSE, bias, and variance of estimates for β_1 at location 5 (**minimum**, *next best*).

Variable name	Description
<code>pag</code>	Proportion working in agriculture
<code>pex</code>	Proportion working in extraction (mining)
<code>pman</code>	Proportion working in manufacturing
<code>pserve</code>	Proportion working in services
<code>pfire</code>	Proportion working in finance, insurance, and real estate
<code>potprof</code>	Proportion working in other professions
<code>pwh</code>	Proportion who are white
<code>pblk</code>	Proportion who are black
<code>phisp</code>	Proportion who are hispanic
<code>metro</code>	Is the county in a metropolitan area?

Table 12: Description of the variables used in the census-data example