

Local Variable Selection and Parameter Estimation of Spatially Varying Coefficient Models

Wesley Brooks

1. Introduction

Whereas the coefficients in traditional linear regression are scalar constants, the coefficients in a varying coefficient regression (VCR) model are functions - often *smooth* functions - of some effect modifying variable (Hastie and Tibshirani, 1993). When the effect modifying variable represents location in a spatial domain, a VCR model implies a spatially local regression model such that the regression coefficients vary over space and will be referred to as a spatially varying coefficient model (SVCR). Statistical inference for the coefficients as functions of location in an SVCR model is more complicated than estimating the coefficients in a global linear regression model where the coefficients are constant across the spatial domain. This document concerns the development of new methodologies for the analysis of spatial data using SVCR.

The methodology described herein is directly applicable only to geostatistical data, which is spatial data observed at discrete locations. Let \mathcal{D} be a spatial domain on which data is collected, and let \mathbf{s} denote a location variable that indexes the domain \mathcal{D} . Let univariate $\{Y(\mathbf{s}) : \mathbf{s} \in \mathcal{D}\}$ and possibly multivariate $\{\mathbf{X}(\mathbf{s}) : \mathbf{s} \in \mathcal{D}\}$ denote random fields called the response and the covariates, respectively. For $i = 1, \dots, n$, let \mathbf{s}_i denote the location in \mathcal{D} of the i th observation of the response and the covariates. Then the data are a realization of the random variables $\{Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n), \mathbf{X}(\mathbf{s}_1), \dots, \mathbf{X}(\mathbf{s}_n)\}$. Let the observed data be denoted $\{y(\mathbf{s}_1), \dots, y(\mathbf{s}_n), \mathbf{x}(\mathbf{s}_1), \dots, \mathbf{x}(\mathbf{s}_n)\}$.

Areal data is a different kind of spatial data in which the spatial domain \mathcal{D} consists of n regions $\{r_1, \dots, r_n\}$. In the case of areal data, the random variables $\{Y(r_1), \dots, Y(r_n), \mathbf{X}(r_1), \dots, \mathbf{X}(r_n)\}$ are defined for regions instead of for points; population and spatial mean temperature are examples of areal data. The analytical method described herein can be applied to areal data if it is recast as geostatistical data by assuming that the data are point-referenced to the centroid of each region, i.e. $\{\mathbf{X}(\mathbf{s}_i), Y(\mathbf{s}_i)\} = \{\mathbf{X}(r_i), Y(r_i)\}$ where \mathbf{s}_i is the centroid of r_i for $i = 1, \dots, n$. The data example in section 5 uses areal data relating to county-level demographics in this way.

Common practice in the analysis of geostatistical and areal data is to model the response variable with a spatial linear regression model consisting of the sum of a fixed mean function, a spatial random effect, and random error all on domain \mathcal{D} , as in:

$$Y(\mathbf{s}) = \mathbf{X}(\mathbf{s})'\boldsymbol{\beta} + W(\mathbf{s}) + \varepsilon(\mathbf{s}) \quad (1)$$

where $\mathbf{X}(\mathbf{s})'\boldsymbol{\beta}$ is the mean function consisting of $\mathbf{X}(\mathbf{s})$, a possibly multivariate spatial random field of covariates, and $\boldsymbol{\beta}$, a vector of regression coefficients. The random error $\varepsilon(\mathbf{s})$ denotes a white noise field such that the errors are independent and identically distributed with mean zero and variance σ^2 , while the random component $W(\mathbf{s})$ denotes a mean-zero, second-order stationary random field that is independent of the random error. The mean function captures the large-scale systematic trend of the response, the spatial random field $W(\mathbf{s})$ can be thought of as a small-scale spatial random effect, and the error term $\varepsilon(\mathbf{s})$ captures micro scale variation (Cressie, 1993). It is common to prespecify the form of a covariance function for the spatial random effect $W(\mathbf{s})$ (Diggle and Ribeiro, 2007). For example, the exponential covariance function (a special case of the Matérn

class of covariance functions) has the form

$$\text{Cov}(W(\mathbf{s}), W(\mathbf{t})) = \exp \left\{ -\phi^{-1} \delta(\mathbf{s}, \mathbf{t}) \right\} \quad (2)$$

where ϕ denotes a range parameter and $\delta(\mathbf{s}, \mathbf{t})$ denotes the Euclidean distance between locations \mathbf{s} and \mathbf{t} . The general form of a Matérn class covariance function is

$$\text{Cov}(W(\mathbf{s}), W(\mathbf{t})) = \left\{ \Gamma(\nu) 2^{\nu-1} \right\}^{-1} \left\{ \delta(\mathbf{s}, \mathbf{t}) \phi^{-1} \sqrt{2\nu} \right\}^{\nu} K_{\nu} \left(\delta(\mathbf{s}, \mathbf{t}) \phi^{-1} \sqrt{2\nu} \right) \quad (3)$$

where ϕ denotes a range parameter, ν denotes the degree of smoothness, K_{ν} denotes the modified Bessel equation of the second kind, and $\delta(\mathbf{s}, \mathbf{t})$ denotes the Euclidean distance between locations \mathbf{s} and \mathbf{t} . The exponential covariance function corresponds to a Matérn class covariance function with $\nu = 1/2$.

A spatial field is said to be stationary if the joint distribution of a sample from the field does not change when the sample locations are all shifted in space by the same amount. Let $\{T(\mathbf{s}) : \mathbf{s} \in \mathcal{D}\}$ be a random field on spatial domain \mathcal{D} that takes value $T(\mathbf{s}_i)$ at location $\mathbf{s}_i \in \mathcal{D}$ for $i = 1, \dots, n$. The random field $T(\mathbf{s})$ is stationary if $F_n(T(\mathbf{s}_1), \dots, T(\mathbf{s}_n)) = F_n(T(\mathbf{s}_1 + \mathbf{h}), \dots, T(\mathbf{s}_n + \mathbf{h}))$ where $F_n(\cdot)$ is the joint distribution of a length n sample from $T(\mathbf{s})$. A spatial random field is second-order stationary if the joint distribution of any two observations from a sample does not change when the sample locations are shifted by the same amount.

The coefficient vector β in (1) is a specific example of the case where $\{\beta(\mathbf{s}) : \mathbf{s} \in \mathcal{D}\}$ is a spatial random field. Specifically, the coefficient vector in (1) represents the case of where $\beta(\mathbf{s}) \equiv \beta$, $\forall \mathbf{s} \in \mathcal{D}$, i.e. the random field is constant. Clearly such random field is stationary. It is also possible to specify a non-constant random coefficient field that is nevertheless stationary. One such model

is written

$$Y(\mathbf{s}) = \mathbf{X}(\mathbf{s})'\boldsymbol{\beta}(\mathbf{s}) + \varepsilon(\mathbf{s}) \quad (4)$$

Where $\boldsymbol{\beta}(\mathbf{s})$ has a Matérn-class covariance function. The random coefficient field $\boldsymbol{\beta}(\mathbf{s})$ can be estimated by Markov Chain Monte Carlo (MCMC) methods (Gelfand et al., 2003).

The preceding methods are limited to the case of a stationary coefficient field.

Both kernel-based and spline-based methods are available for fitting varying coefficient models. For example, it is straightforward to modify a thin plate regression spline model into a VCR model (Wood, 2006). Alternatively, the local likelihood can be used to fit generalized linear models with varying coefficients using kernel smoothing (Loader, 1999). Fan and Zhang (1999) demonstrated that the optimal kernel bandwidth estimate for a VCR model can be found via a two-step technique.

Model selection in VCR models may be local or global. Global selection means including or excluding variables everywhere in the model domain, while local selection means including or excluding variables at each observation location. Two methods have been proposed for global model selection in spline-based VCR models. Wang et al. (2008) applied a SCAD penalty (Fan and Li, 2001) for variable selection in spline-based VCR models with a univariate effect-modifying variable. Antoniadis et al. (2012) used the nonnegative Garrote penalty (Breiman, 1995) in P-spline-based VCR models having a univariate effect-modifying variable.

Wavelet methods for fitting SVCR models were explored by ? and Zhang and Clayton (2011). Sparsity in the wavelet coefficients is achieved either by ℓ_1 -penalization (?) or by Bayesian variable selection (Zhang and Clayton, 2011). Sparsity in the wavelet domain does not imply sparsity in

the covariates, though, so neither method can be used for local variable selection.

Geographically Weighted Regression (GWR) is a kernel-based method of estimating the coefficients of a VCR model in the context of spatial data (Brundson et al., 1998; Fotheringham et al., 2002). GWR uses kernel-weighted regression with weights based on the distance between observation locations. The presentation of GWR in Fotheringham et al. (2002) followed the development of local likelihood in Loader (1999). GWR can be thought of as a kernel smoother for regression coefficients, which tends to exhibit bias near the boundary of the region being modeled (Hastie and Loader, 1993). One way to reduce the boundary-effect bias is to model the coefficient surface as locally linear rather than locally constant by including coefficient-by-location interactions (Hastie and Loader, 1993). Adding these interactions to the GWR model is analogous to a transition from kernel smoothing to local regression, which was introduced in Wang et al. (2008).

GWR relies on *a priori* global model selection to decide which variables should be included in the model. In the context of ordinary least squares regression, ℓ_1 regularization for variable selection is called the Lasso (Tibshirani, 1996). While popular, the Lasso does not generally produce consistent estimates of the relevant predictor variables (Leng et al., 2006). Regularization methods such as the Adaptive Lasso (AL) (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 idea of using ℓ_1 regularization for local variable selection in a GWR model has appeared in the literature as the Geographically Weighted Lasso (GWL) (Wheeler, 2009). The GWL uses the Lasso with a jackknife criterion for selection of the 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.

This paper introduces a method of regularization for local variable selection in GWR models that avoids this limitation of the GWL by using a penalized-likelihood criterion to select the Lasso tuning parameters. Here we use a version of the BIC, but in principle one could use another information criterion like the AIC. The local BIC presented here is based on the local likelihood (Loader, 1999) and the total BIC is based on an *ad hoc* calculation of the sample size and degrees of freedom for estimating the spatially-varying coefficient surfaces.

Three regularization methods were used in this work. The AL was implemented in two ways - once via the lars algorithm (Efron et al., 2004) which uses least squares, and once via coordinate descent using the R package glmnet (Friedman et al., 2010). The third regularization method implemented here uses the Adaptive Elastic Net (AEN) penalty (Zou and Zhang, 2009), also via coordinate descent using the glmnet package.

2. Geographically Weighted Regression

2.1. Model

Consider n data observations, taken 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 $\boldsymbol{\beta}(\mathbf{s}_i)$ may be spatially-varying and $\varepsilon(\mathbf{s}_i)$ is random noise at location \mathbf{s}_i . That is,

$$y(\mathbf{s}_i) = \mathbf{x}(\mathbf{s}_i)' \boldsymbol{\beta}(\mathbf{s}_i) + \varepsilon(\mathbf{s}_i) \tag{5}$$

Further assume that the error term $\varepsilon(\mathbf{s}_i)$ is normally distributed with zero mean and variance σ^2 ,

and that $\varepsilon(\mathbf{s}_i)$, $i = 1, \dots, n$ are independent.

$$\varepsilon(\mathbf{s}_i) \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2) \quad (6)$$

In order to simplify the notation, let $\mathbf{x}(\mathbf{s}_i) \equiv \mathbf{x}_i \equiv (1, x_{i1}, \dots, x_{ip})'$, $\boldsymbol{\beta}(\mathbf{s}_i) \equiv \boldsymbol{\beta}_i \equiv (\beta_{i0}, \beta_{i1}, \dots, \beta_{ip})'$, and $y(\mathbf{s}_i) \equiv y_i$. Further, let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$ and $\mathbf{y} = (y_1, \dots, y_n)'$. Equations (5) and (6) can now be rewritten as

$$y_i = \mathbf{x}_i' \boldsymbol{\beta}_i + \varepsilon_i \text{ and } \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2) \quad (7)$$

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

$$\ell(\boldsymbol{\beta}) = -(1/2) \left\{ n \log(2\pi\sigma^2) + (\sigma^2)^{-1} \sum_{i=1}^n (y_i - \mathbf{x}_i' \boldsymbol{\beta}_i)^2 \right\} \quad (8)$$

Since there are a total of $n \times (p+1)$ free parameters for n observations 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 coefficients $\boldsymbol{\beta}(\mathbf{s})$ are smoothly varying over space, 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 the traditional GWR.

2.2. Estimation

In geographically weighted regression, the coefficient surface $\boldsymbol{\beta}(\mathbf{s})$ is estimated at each sampling 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' . The bi-square kernel can be used to generate spatial weights based on the Euclidean distances and a bandwidth ϕ . The bisquare kernel assigns the maximum weight of one where $\mathbf{s}_i = \mathbf{s}_{i'}$ so $\delta_{ii'} = 0$, discontinuously differentiable, and assigns zero weight to observations at distances greater than one bandwidth from \mathbf{s}_i :

$$w_{ii'} = \begin{cases} \left[1 - (\phi^{-1}\delta_{ii'})^2\right]^2 & \text{if } \delta_{ii'} < \phi \\ 0 & \text{if } \delta_{ii'} \geq \phi \end{cases} \quad (9)$$

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

$$\mathcal{L}_i(\boldsymbol{\beta}_i) = \prod_{i'=1}^n \left[(2\pi\sigma_i^2)^{-1/2} \exp \left\{ - (2\sigma_i^2)^{-1} (y_{i'} - \mathbf{x}_{i'}'\boldsymbol{\beta}_i)^2 \right\} \right]^{w_{ii'}} \quad (10)$$

where σ_i^2 is a local approximation to the error variance σ^2 . Thus, the local log-likelihood function is:

$$\ell_i(\boldsymbol{\beta}_i) \propto - (1/2) \sum_{i'=1}^n w_{ii'} \left\{ \log \sigma_i^2 + (\sigma_i^2)^{-1} (y_{i'} - \mathbf{x}_{i'}'\boldsymbol{\beta}_i)^2 \right\} \quad (11)$$

The GWR coefficient estimates $\hat{\boldsymbol{\beta}}_{i,\text{GWR}}$ maximize the local likelihood at location \mathbf{s}_i . From (10) and (11), it is apparent that $\hat{\boldsymbol{\beta}}_{i,\text{GWR}}$ can be calculated using weighted least squares. Let \mathbf{W}_i denote a diagonal weight matrix with

$$\mathbf{W}_i = \text{diag} \{w_{ii'}\}_{i'=1}^n \quad (12)$$

Thus, it follows that

$$\hat{\beta}_{i,\text{GWR}} = (\mathbf{X}'\mathbf{W}_i\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}_i\mathbf{y} \quad (13)$$

The estimate of σ_i^2 is attained by maximizing (11). Thus,

$$\begin{aligned} \hat{\sigma}_i^2 &= (\mathbf{1}'_n \mathbf{w}_i)^{-1} \left(\mathbf{y} - \mathbf{X} (\mathbf{X}'\mathbf{W}_i\mathbf{X})^{-1} \mathbf{X}'\mathbf{W}_i\mathbf{y} \right)' \mathbf{W}_i \left(\mathbf{y} - \mathbf{X} (\mathbf{X}'\mathbf{W}_i\mathbf{X})^{-1} \mathbf{X}'\mathbf{W}_i\mathbf{y} \right) \\ &= (\mathbf{1}'_n \mathbf{w}_i)^{-1} (\mathbf{y} - \hat{\mathbf{y}})' \mathbf{W}_i (\mathbf{y} - \hat{\mathbf{y}}) \end{aligned} \quad (14)$$

where $\mathbf{1}_n$ is a length n column vector of ones.

3. Model Selection

3.1. Local Variable Selection

Both the AL and the AEN are explored for local variable selection in GWR models. In fact, since the AL is an ℓ_1 regularization method while the AEN is a combine ℓ_1 and ℓ_2 regularization method, the AL is a special case of the AEN where the ℓ_2 penalty is set to zero.

The AL is an ℓ_1 regularization method for variable selection in regression models (Zou, 2006). Unlike the traditional Lasso (Tibshirani, 1996), which applies an equal penalty λ_i^* to each covariate in the local model at \mathbf{s}_i , the AL adjusts the penalty of each covariate based on the covariate's unpenalized local coefficient.

The AEN generalizes the AL penalty to include an additional ridge penalty (Zou and Zhang, 2009). Ridge regression is an ℓ_2 regularization technique that differs from the Lasso in that the ridge penalty λ_i^\dagger is applied to the sum of the squared local regression coefficients (Hoerl and Kennard, 1970). The ridge penalty is used to estimate coefficients in regression models with correlated co-

variates because it stabilizes the inversion of the covariance matrix, which robustifies the coefficient estimates (Hastie et al., 2009).

3.1.1. Adaptive Lasso

The objective minimized to fit the local Geographically Weighted Adaptive Lasso (GWAL) model at \mathbf{s}_i is

$$\mathcal{S}_i = \sum_{i'=1}^n w_{ii'} (y_{i'} - \mathbf{x}_{i'}' \boldsymbol{\beta}_i)^2 + \lambda_i \sum_{j=1}^p |\beta_{ij}/\gamma_{ij}| \quad (15)$$

where $\sum_{i'=1}^n w_{ii'} (y_{i'} - \mathbf{x}_{i'}' \boldsymbol{\beta}_i)^2$ is the weighted least squares objective minimized by traditional GWR, and $\lambda_i \sum_{j=1}^p |\beta_{ij}/\gamma_{ij}|$ is the AL penalty. Letting the vector of unpenalized local coefficients be $\boldsymbol{\gamma}_i$, the AL penalty for covariate j at location \mathbf{s}_i is $|\lambda_i/\gamma_{ij}|$, where λ_i is a the local penalty that applies to all coefficients at location \mathbf{s}_i and $\boldsymbol{\gamma}_i = \{|\gamma_{ij}|\}$ is the vector of adaptive weights at location \mathbf{s}_i .

3.1.2. Adaptive Elastic Net

The objective minimized to fit the local Geographically Weighted Adaptive Elastic Net (GWAEN) model at \mathbf{s}_i is

$$\begin{aligned} \mathcal{S}_i &= \sum_{i'=1}^n w_{ii'} (y_{i'} - \mathbf{x}_{i'}' \boldsymbol{\beta}_i)^2 + \alpha_i \lambda_i \sum_{j=1}^p |\beta_{ij}/\gamma_{ij}| + (1 - \alpha_i) \lambda_i \sum_{j=1}^p (\beta_{ij}/\gamma_{ij})^2 \\ &= \sum_{i'=1}^n w_{ii'} (y_{i'} - \mathbf{x}_{i'}' \boldsymbol{\beta}_i)^2 + \lambda_i \left(\alpha_i \sum_{j=1}^p |\beta_{ij}/\gamma_{ij}| + (1 - \alpha_i) \sum_{j=1}^p [\beta_{ij}/\gamma_{ij}]^2 \right) \end{aligned} \quad (16)$$

where the adaptive weights $\boldsymbol{\gamma}_i = \{|\gamma_{ij}|\}$ are calculated as for the AL, and the elastic net parameter α_i controls the balance between the ℓ_1 and ℓ_2 penalties.

Fitting a GWAEN model requires selecting the vector of elastic net parameters $\boldsymbol{\alpha} = \{\alpha_i\}$. In the simulation study (Section 4), the elastic net parameter is chosen globally ($\alpha_i \equiv \alpha$ for $i = 1, \dots, n$).

The global elastic net parameter is calculated from the maximum global (i.e. for all data without weighting) Pearson correlation between any two covariates, ρ_{\max} : $\alpha = 1 - \rho_{\max}$.

3.2. Tuning Parameter Selection

Each local model in a GWAL or GWAEN model has a local tuning parameter λ_i . To select λ_i , we propose a locally-weighted version of the Bayesian Information Criterion (BIC) (Schwarz, 1978) which we call the local BIC (BIC_{loc}):

$$\begin{aligned}
\text{BIC}_{\text{loc},i} &= -2 \sum_{i'=1}^n \ell_{ii'} + \left(\sum_{i'=1}^n w_{ii'} \right) \text{df}_i \\
&= -2 \times \sum_{i'=1}^n \log \left\{ (2\pi \hat{\sigma}_i^2)^{-1/2} \exp \left[-\frac{1}{2} \hat{\sigma}_i^{-2} (y_{i'} - \mathbf{x}'_{i'} \hat{\boldsymbol{\beta}}_i)^2 \right] \right\}^{w_{ii'}} + \left(\sum_{i'=1}^n w_{ii'} \right) \text{df}_i \\
&= \sum_{i'=1}^n w_{ii'} \left\{ \log(2\pi) + \log \hat{\sigma}_i^2 + \hat{\sigma}_i^{-2} (y_{i'} - \mathbf{x}'_{i'} \hat{\boldsymbol{\beta}}_i)^2 \right\} + \left(\sum_{i'=1}^n w_{ii'} \right) \text{df}_i \\
&= \hat{\sigma}_i^{-2} \sum_{i'=1}^n w_{ii'} (y_{i'} - \mathbf{x}'_{i'} \hat{\boldsymbol{\beta}}_i)^2 + \left(\sum_{i'=1}^n w_{ii'} \right) \text{df}_i + C_i
\end{aligned} \tag{17}$$

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

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 BIC 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. Coefficient Estimation

Having selected the variables for inclusion in the model, either by the GWAL or the GWAEN, the model's coefficient estimates are computed via weighted least squares on the selected variables without regularization. That is, the local coefficient estimates are:

$$\hat{\beta}_i = \underset{\beta}{\operatorname{argmin}} \sum_{i'=1}^n w_{ii'} (y_{i'} - \mathbf{x}_{i'}' \boldsymbol{\Omega}_i \beta)' \quad (18)$$

$$= \left(\tilde{\mathbf{X}}' \mathbf{W}_i \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}' \mathbf{W}_i \mathbf{y} \quad (19)$$

where $\boldsymbol{\Omega}_i = \operatorname{diag}(\boldsymbol{\omega}_i)$; $\boldsymbol{\omega}_i = \{\omega_{ij}\}$ is the active set vector of indicators $\omega_{ij} = \mathbf{I}(\beta_{ij} \neq 0)$; and $\tilde{\mathbf{X}}_i$ is the design matrix \mathbf{X} with columns corresponding to zeroes in $\boldsymbol{\omega}_i$ removed.

3.4. Bandwidth selection

Letting $H_i = \left\{ \mathbf{W}_i^{1/2} \tilde{\mathbf{X}} \left(\tilde{\mathbf{X}}' \mathbf{W}_i \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}' \mathbf{W}_i^{1/2} \right\}_i$ where $\{\cdots\}_i$ represents the i th row of the enclosed matrix, we have

$$\mathbf{H} = (H_1 \cdots H_n)' \quad (20)$$

and the fitted values from the model are

$$\hat{\mathbf{Y}} = \mathbf{H} \mathbf{Y} \quad (21)$$

The global bandwidth parameter ϕ in (9) is selected by minimizing an approximation to the global AIC:

$$\text{AIC} = 2n \log \sigma + n \left\{ \frac{n + \nu}{n - 2 - \nu} \right\} \quad (22)$$

where ν is the trace of the smoothing matrix H , and approximates the total degrees of freedom of the GWAL or GWAEN model. (Hurvich et al., 1998).

4. Simulation

4.1. Simulation Setup

A simulation study was conducted to assess the performance 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 X_1, \dots, X_5 was simulated by a Gaussian random field (GRF) with mean zero and exponential spatial covariance $\text{Cov}(X_{ji}, X_{ji'}) = \sigma_x^2 \exp(-\tau_x^{-1} \delta_{ii'})$ where $\sigma_x^2 = 1$ is the variance, $\tau_x = 0$ 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{X} 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 $y_i = \mathbf{x}_i' \boldsymbol{\beta}_i + \varepsilon_i$ for $i = 1, \dots, n$ where $n = 900$ and for simplicity the ε_i 's were iid Gaussian with mean zero and variance σ_ε^2 .

The simulated data include the output y and five covariates X_1, \dots, X_5 . The true data-generating

model uses only X_1 , so X_2, \dots, X_5 are included to assess performance in variable-selection.

There were twelve simulation settings, each of which was simulated 100 times. For each of the twelve settings, $\beta_1(\mathbf{s})$, the true coefficient surface for \mathbf{X}_1 , was nonzero in at least part of the simulation domain. There were four other simulated covariates, but their true coefficient surfaces were zero across the area under simulation. The twelve simulation settings are described in Table 1. Three parameters were varied to produce the twelve settings: there were three functional forms for the coefficient surface $\beta_1(\mathbf{s})$ (step, gradient, and parabola - see Figure 1); data was simulated both with ($\rho = 0.5$) and without ($\rho = 0$) correlation between the covariates; and simulations were made with low ($\sigma_\varepsilon^2 = 0.25$) and high ($\sigma_\varepsilon^2 = 1$) variance for the random noise term.

The performance of the penalized GWR methods (AL via `lars` and via `glmnet`, and the AEN (`enet`) was compared to that of 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. Also included in the comparison was the GWR algorithm of Fotheringham et al. (2002) without variable selection (`gwr`). Finally, there is a category of simulation results using the three penalized GWR methods for local variable selection and then ordinary GWR for coefficient estimation.

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

4.2. Results

Selection. Table 2 lists the results of variable selection. The correct variable was usually included in the local models, and the unimportant variables were usually excluded. Arguably the least-accurate selection was at locations one and five for the step function using the **lars** algorithm, where variables that do not appear in the true model were selected for inclusion at rates between 11% and 22%. The **enet** and **glmnet** algorithms, using the same data, had false-positive errors at rates between 0% and 8%, which are typical of the error rates for all other location/function/algorithm combinations.

Selection performance was more affected by an increase in the noise variance from $\sigma_\varepsilon = 0.5$ to $\sigma_\varepsilon = 1$ than by an increase in collinearity from $\rho = 0$ to $\rho = 1$. For instance, for the step function at location three, $\beta_1(\mathbf{s}_3) = 0.5$. Where $\sigma_\varepsilon = 0.5$, the **glmnet** algorithm selected $\beta_1(\mathbf{s}_3)$ for inclusion at a rate of 100% (when $\rho = 0$) and 99% (when $\rho = 0.5$). But when $\sigma_\varepsilon = 1$, the rate of selection for $\beta_1(\mathbf{s}_3)$ at a rate of 75% (when $\rho = 0$) and 68% (when $\rho = 0.5$).

The **enet** algorithm outperforms the others in selection but the difference is small - a roughly one percentage point improvement in the rate of true positives and true negatives when $\rho = 0.5$. There is no apparent difference between **glmnet** and **enet** when $\rho = 0$.

Coefficient Estimation. The MSE, bias, and variance of $\hat{\beta}_1$ are listed in Tables 3, 4, and 5, respectively. The method of oracular selection led to the best MSE in 28 of the 60 cases, which is more than any other single method. In general, the methods that do local variable selection had lower MSE than the basic GWR. As was the case for selection, estimation accuracy (in terms

of MSE) suffered more by an increase in σ_ϵ from 0.5 to 1 than from an increase in ρ from 0 to 0.5. Oracular selection was decisively superior to base GWR and to local variable selection for estimating the gradient β_1 , turning in the best MSE for all combinations of location and simulation parameters.

In general, oracular selection and base GWR were quite similar in terms of $\text{var}(\hat{\beta}_1)$, with notably greater variance for the local selection methods. However, the local selection methods had less bias than base GWR, even exhibiting less bias than oracular selection in many settings. There was no simulation setting for which base GWR had the smallest or second-smallest bias.

It seems, therefore, that the local selection methods reduce bias and increase variance of the coefficient estimates, as compared to base GWR. Whether base GWR or local selection is better in terms of MSE of the coefficient estimates is not clear in all cases, but when the actual coefficient is equal to zero (or nearly so), local selection does seem to reduce the MSE over base GWR.

Fitted Values. The MSE of the \hat{Y} , $\text{MSE}(\hat{Y})$, is listed in Table 6. Nominally, $\text{MSE}(\hat{Y})$ should be equal to the noise variance, σ_ϵ^2 , which is 1 for odd-numbered rows and 0.25 for even numbered rows. There is not much difference in $\text{MSE}(\hat{Y})$ between the various estimation methods, except that it is larger for the oracular and **gwr** methods where $\beta_1(\mathbf{s})$ is near or equal to zero.

4.3. Tables

4.3.1. Selection

4.3.2. Estimation

5. Data Analysis

5.1. Census Poverty Data

An example data analysis is presented to demonstrate application of penalized GWR. In this example we use penalized GWR to do local variable selection and coefficient estimation for a varying-coefficients model of how poverty is related to a list of demographic and social variables. The data is from the U.S. Census Bureau's decennial census in the year 1970. This analysis looks specifically at the upper midwestern states of Minnesota, Iowa, Wisconsin, Illinois, Indiana, and Michigan. This is areal data, aggregated at the county level.

Three kinds of variables were considered as potential predictors of county-level poverty rate.

- Variables that describe the county's employment structure (**pag**, the proportion of residents employed in agriculture, **pex**, the proportion of residents employed in mining, **man**, the proportion of residents employed in manufacturing, **pfire**, the proportion of residents employed in finance, insurance, and real estate, **pserve**, the proportion of residents employed in services, and **potprof**, the proportion of residents employed in other professions)
- Variables that describe the county's racial makeup (**pwh**, the proportion of residents who are white, **pblk**, the proportion of residents who are black, and **phispanic**, the proportion of residents who are hispanic)

- `pmetro`: an indicator of whether the county is in a metropolitan area.

The outcome of interest (poverty rate) is a proportion, taking values in $[0, 1]$. To demonstrate the geographically-weighted Lasso in a linear regression context, we model the logit-transformed poverty rate. The predictor variables were not transformed - raw proportions were used.

5.2. Modeling

The AEN was used for variable selection, and then coefficients for the selected variables were estimated by weighted least squares without shrinkage. The standard `gwr` algorithm was used to fit a model to the same data for the sake of comparison.

5.3. Figures

The coefficient estimates are plotted on maps of the upper midwest in Figure 3 (based on the AEN) and Figure 4 (for standard GWR).

5.4. Discussion

It is immediately apparent that the estimated coefficient surfaces are non-constant for most variables. The same large-scale patterns appear in both figures, but with differences. First of all, the AEN has selected a larger bandwidth than base GWR, so there is less variability in the coefficient estimates from the AEN. This may be one reason that the AEN coefficient estimates are less extreme than those for base GWR. In a model with a logit-transformed proportion as the output, the coefficients can be interpreted as log odds ratios, so, e.g., the estimate of -100 as the coefficient of `phisp` (albeit at the edge of the domain) seems unrealistic.

Assessing variable selection for this data is difficult, since the AEN almost never removed any variables from the model. Indeed, some coefficients seem nearly constant across the domain. An exception is the coefficient surface for `pex` (mining employment). That surface indicates an interaction whereby the proportion of people working in mining in southern parts of the domain is associated with an increase in the poverty rate, while in northern parts of the domain it is associated with a decrease in the poverty rate.

6. References

References

- 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.
- Brundson, C., S. Fotheringham, and M. Charlton (1998). Geographically weighted regression: a natural evolution of the expansion method for spatial data analysis. *Environment and Planning A* 30, 1905–1927.
- Cressie, N. (1993). *Statistics for spatial data*. Wiley.
- Diggle, P. and P. Ribeiro (2007). *Model-based geostatistics*. Springer New York.
- Efron, B., T. Hastie, I. Johnstone, and R. Tibshirani (2004). Least angle regression. *Annals of Statistics* 32(2), 407–499.

- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* 96(456), 1348–1360.
- 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.
- Friedman, J., T. Hastie, and R. Tibshirani (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* 33(1), 1–22.
- Gelfand, A. E., H.-J. Kim, C. F. Sirmans, and S. Banerjee (2003). Spatial modeling with spatially varying coefficient processes. *Journal of the American Statistical Association* 98(462), 387–396.
- 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.
- Hastie, T., R. Tibshirani, and J. Friedman (2009). *The elements of statistical learning: data mining, inference, and prediction*. Springer New York.
- Hoerl, A. E. and R. W. Kennard (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* 12(1), 55–67.
- Hurvich, C. M., J. S. Simonoff, and C.-L. Tsai (1998). Smoothing parameter selection in nonparametric regression using an improved akaike information criterion. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 60(2), pp. 271–293.

- Leng, C., Y. Lin, and G. Wahba (2006). A note on the lasso and related procedures in model selection. *Statistica Sinica* 16, 1273–1284.
- Loader, C. (1999). *Local regression and likelihood*. Springer New York.
- Schwarz, G. (1978). Estimating the dimensions of a model. *The Annals of Statistics* 6(2), 461–464.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* 58, 267–288.
- 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.
- Zhang, J. and M. Clayton (2011). Functional concurrent linear regression model for images. *Journal of Agricultural, Biological, and Environmental Statistics* 16(1), 105–130.
- 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.

Zou, H. and H. Zhang (2009). On the adaptive elastic net with a diverging number of parameters. *The Annals of Statistics* 37(4), 1733–1751.

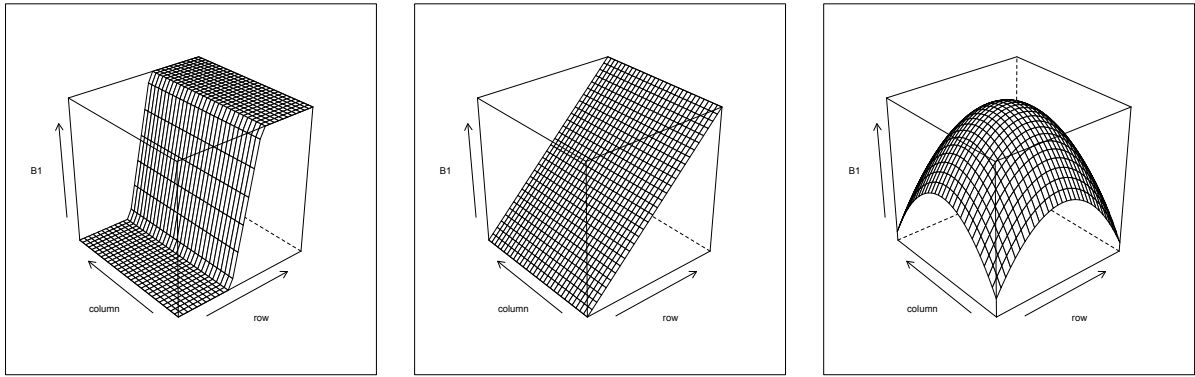


Figure 1: The actual β_1 coefficient surface used in the simulation.

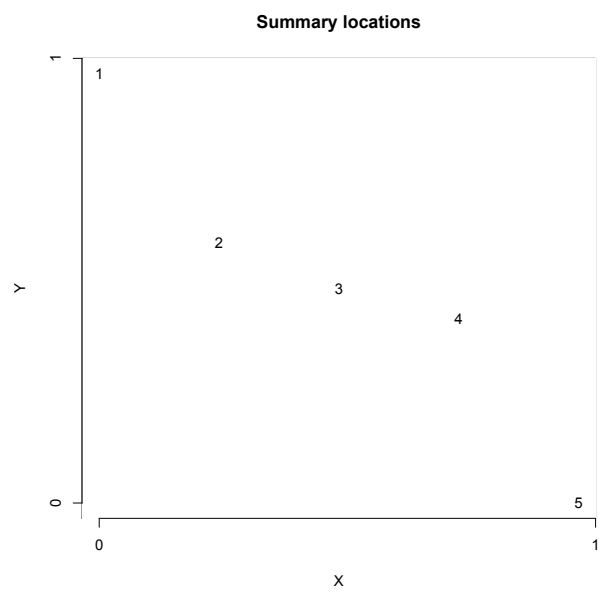


Figure 2: Locations where the variable selection and coefficient estimation of GWL were summarized.

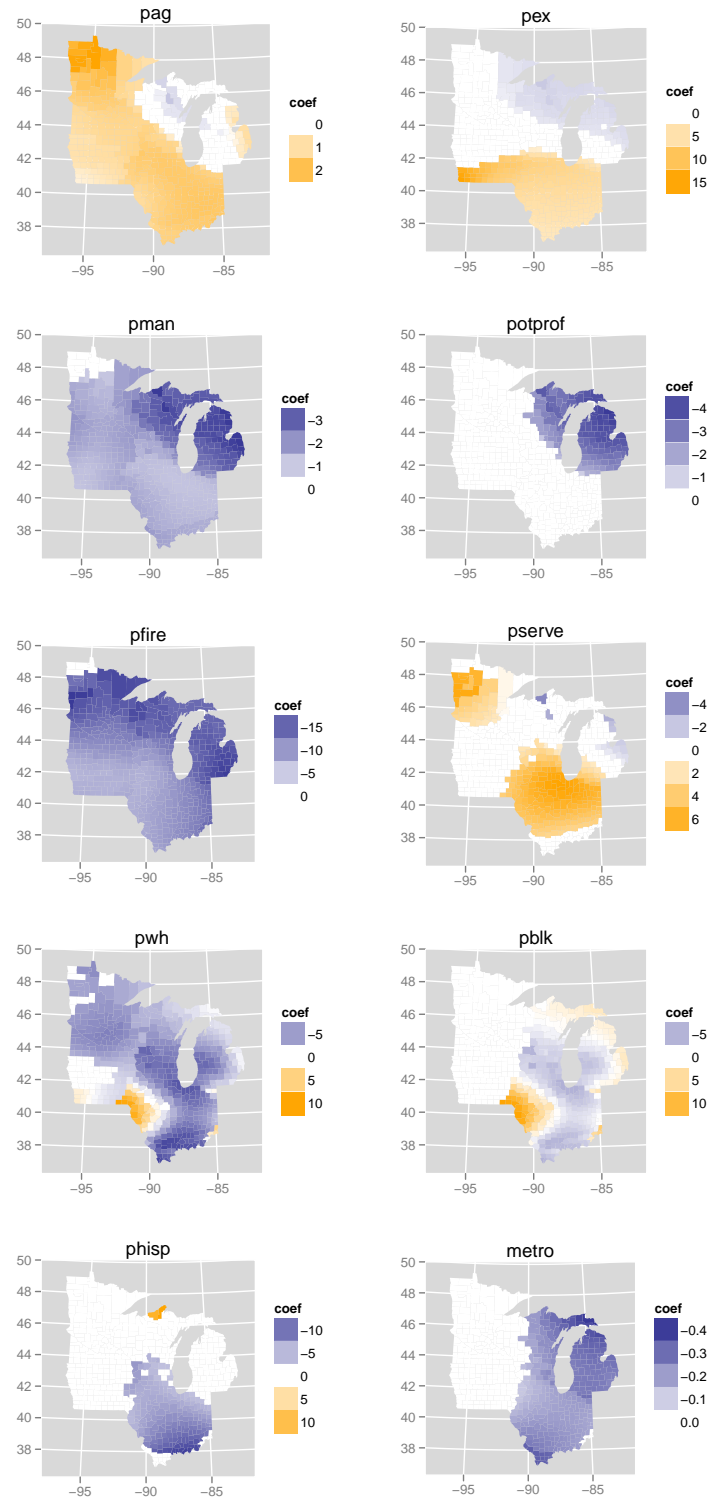


Figure 3: Coefficient surfaces for the logit of poverty rate, based on the 1970 census and estimated by the unshrunk AEN.

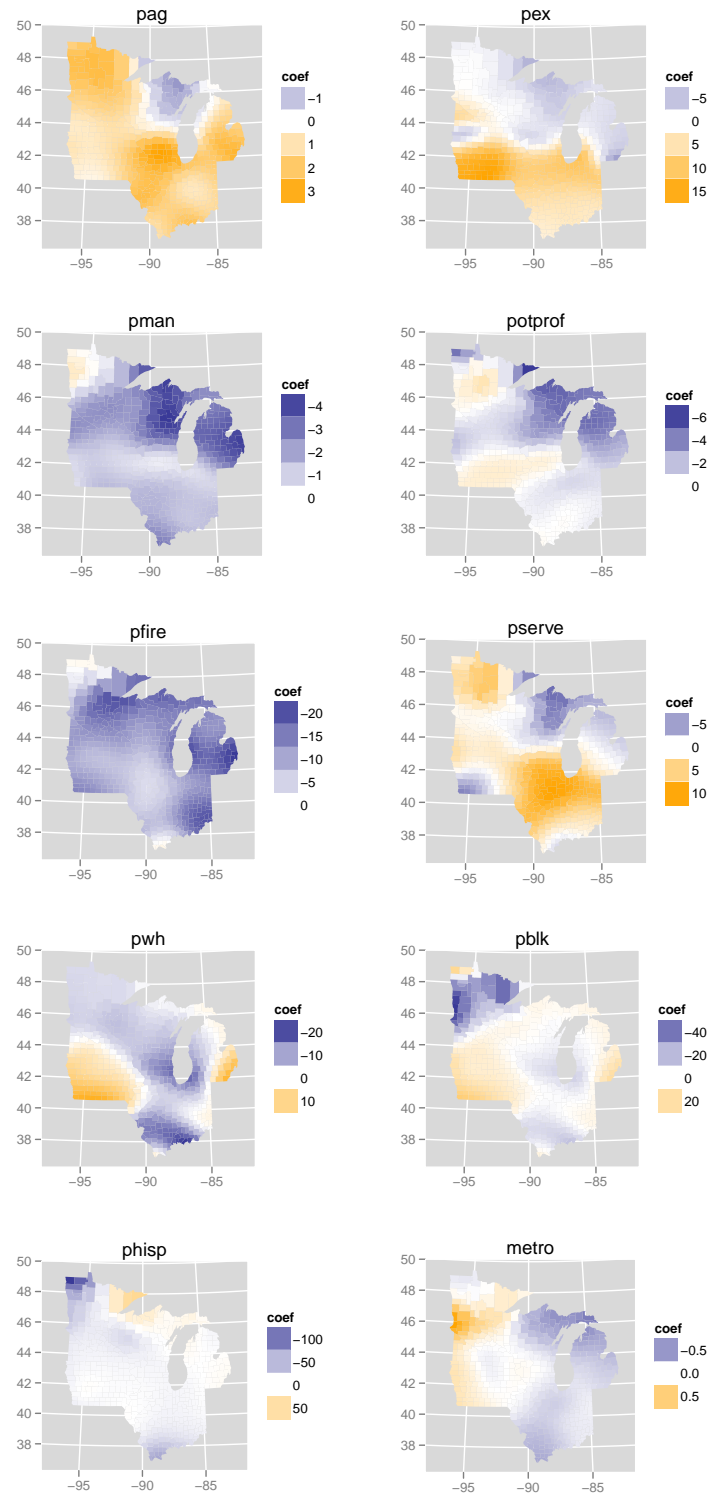


Figure 4: Coefficient surfaces for the logit of poverty rate based on the 1970 census and estimated by base GWR.

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 1: Simulation parameters for each setting.

location	step						gradient						parabola					
	lars			enet			glmnet			lars			enet			lars		
	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	1.00	0.14		1.00	0.00		1.00	0.00		1.00	0.02		1.00	0.00		0.73	0.08	
	1.00	0.16		0.90	0.06		0.90	0.06		1.00	0.02		1.00	0.02		0.72	0.05	
	1.00	0.19		1.00	0.00		1.00	0.00		1.00	0.02		1.00	0.00		0.73	0.08	
	0.99	0.11		0.91	0.08		0.90	0.06		0.99	0.04		0.97	0.04		0.62	0.08	
2	1.00	0.06		1.00	0.00		1.00	0.00		1.00	0.01		1.00	0.00		1.00	0.03	
	1.00	0.06		0.97	0.03		0.97	0.02		1.00	0.02		1.00	0.01		0.99	0.02	
	1.00	0.04		1.00	0.00		1.00	0.00		1.00	0.02		1.00	0.00		1.00	0.02	
	1.00	0.01		0.98	0.04		0.97	0.02		1.00	0.03		0.96	0.02		0.95	0.03	
3	1.00	0.04		1.00	0.00		1.00	0.00		1.00	0.02		1.00	0.00		1.00	0.01	
	0.94	0.05		0.74	0.02		0.75	0.02		1.00	0.01		0.98	0.00		1.00	0.02	
	0.99	0.04		0.98	0.00		0.99	0.00		1.00	0.01		1.00	0.00		1.00	0.01	
	0.87	0.05		0.68	0.04		0.68	0.04		0.97	0.01		0.91	0.02		0.97	0.04	
4	0.72	0.04		0.39	0.00		0.40	0.00		1.00	0.01		1.00	0.00		1.00	0.02	
	0.47	0.03		0.37	0.01		0.38	0.01		1.00	0.01		0.97	0.01		0.98	0.03	
	0.58	0.07		0.34	0.00		0.35	0.00		1.00	0.01		1.00	0.00		1.00	0.01	
	0.51	0.06		0.43	0.03		0.37	0.04		0.97	0.02		0.93	0.02		0.96	0.05	
5	0.17	0.14		0.00	0.00		0.00	0.00		0.93	0.03		0.72	0.00		0.75	0.10	
	0.15	0.14		0.05	0.04		0.05	0.06		0.72	0.05		0.62	0.01		0.73	0.04	
	0.14	0.22		0.00	0.00		0.00	0.00		0.89	0.04		0.76	0.00		0.73	0.09	
	0.14	0.11		0.03	0.03		0.05	0.04		0.71	0.05		0.59	0.02		0.61	0.08	

Table 2: Selection frequency for the simulation experiment

function	location	lars	enet	glmnet	u.lars	u.enet	u.glmnet	oracular	gwr
step	1	0.016	<i>0.009</i>	0.009	0.077	0.055	0.055	0.059	0.111
		0.048	0.150	0.150	0.267	0.566	0.584	<i>0.063</i>	0.110
		0.025	<i>0.021</i>	0.012	0.134	0.072	0.068	0.054	0.114
		<i>0.078</i>	0.182	0.171	0.294	0.434	0.439	0.064	0.112
	2	<i>0.025</i>	0.036	0.036	0.021	0.032	0.032	0.092	0.166
		<i>0.056</i>	0.085	0.082	0.044	0.069	0.069	0.093	0.165
		<i>0.024</i>	0.040	0.034	0.019	0.029	0.028	0.094	0.169
		<i>0.065</i>	0.095	0.101	0.048	0.071	0.081	0.091	0.164
	3	0.007	0.005	0.005	0.004	<i>0.003</i>	0.003	0.004	0.007
		0.021	0.056	0.055	0.023	0.064	0.063	0.005	<i>0.009</i>
		0.009	0.008	0.006	0.007	0.007	<i>0.005</i>	0.004	0.007
		0.033	0.064	0.066	0.033	0.063	0.069	0.006	<i>0.010</i>
	4	0.016	<i>0.019</i>	0.020	0.021	0.022	0.022	0.091	0.166
		0.035	<i>0.033</i>	0.032	0.041	0.038	0.037	0.094	0.169
		0.014	<i>0.016</i>	0.019	0.018	0.020	0.020	0.095	0.166
		0.035	0.027	<i>0.032</i>	0.042	0.039	0.034	0.094	0.171
	5	0.005	0.000	0.000	0.031	0.000	0.000	0.000	0.110
		0.021	0.012	<i>0.012</i>	0.104	0.064	0.061	0.000	0.112
		0.007	0.000	0.000	0.049	0.000	0.000	0.000	0.112
		0.020	<i>0.008</i>	0.012	0.129	0.030	0.040	0.000	0.117
gradient	1	0.113	0.090	0.090	0.007	<i>0.006</i>	0.006	0.003	0.147
		0.103	0.100	0.101	<i>0.045</i>	0.083	0.088	0.010	0.146
		0.108	0.099	0.100	0.013	0.004	<i>0.004</i>	0.002	0.146
		0.117	0.138	0.143	<i>0.079</i>	0.165	0.136	0.009	0.153
	2	0.005	0.003	0.003	0.001	<i>0.001</i>	0.001	0.001	0.005
		0.005	0.005	0.005	<i>0.003</i>	0.003	0.003	0.002	0.005
		0.004	0.004	0.004	0.001	0.001	<i>0.001</i>	0.001	0.005
		0.013	0.022	0.027	0.008	0.019	0.024	0.002	<i>0.007</i>
	3	0.001	0.001	0.001	0.001	0.000	<i>0.000</i>	0.000	0.001
		0.003	0.009	0.009	0.002	0.011	0.011	0.001	<i>0.002</i>
		0.001	0.001	0.000	0.001	0.000	<i>0.000</i>	0.000	0.001
		0.013	0.024	0.021	0.011	0.024	0.020	0.001	<i>0.002</i>
	4	0.004	0.003	0.003	0.001	<i>0.001</i>	0.001	0.001	0.005
		0.007	0.010	0.010	<i>0.003</i>	0.007	0.007	0.002	0.007
		0.004	0.003	0.003	0.002	0.001	<i>0.001</i>	0.001	0.006
		0.013	0.018	0.020	0.010	0.018	0.018	0.002	<i>0.007</i>
	5	0.112	0.079	0.079	0.005	<i>0.003</i>	0.003	0.000	0.151
		0.091	0.071	0.073	0.018	<i>0.012</i>	0.012	0.000	0.152
		0.102	0.079	0.087	0.006	<i>0.004</i>	0.004	0.000	0.148
		0.094	0.056	0.069	0.067	<i>0.012</i>	0.015	0.000	0.153
parabola	1	0.033	0.019	0.019	0.026	0.009	<i>0.009</i>	0.059	0.112
		0.076	0.057	0.057	0.142	0.047	<i>0.049</i>	0.067	0.114
		0.032	0.015	0.020	0.032	0.008	<i>0.008</i>	0.055	0.107
		0.079	0.035	0.041	0.180	<i>0.039</i>	0.042	0.075	0.115
	2	0.004	<i>0.003</i>	0.003	0.004	0.003	0.003	0.010	0.009
		0.014	0.018	0.018	0.015	0.023	0.022	<i>0.012</i>	0.010
		0.004	0.005	0.004	0.005	0.004	<i>0.004</i>	0.012	0.011
		0.022	0.017	0.022	0.024	0.015	0.025	<i>0.011</i>	0.011
	3	0.005	0.005	0.005	0.005	<i>0.004</i>	0.004	0.027	0.027
		<i>0.022</i>	0.025	0.024	0.021	0.024	0.023	0.027	0.028
		0.007	0.007	0.006	0.006	<i>0.005</i>	0.005	0.028	0.030
		0.029	0.021	<i>0.018</i>	0.027	0.015	0.022	0.027	0.029
	4	0.004	0.004	0.004	0.004	0.004	<i>0.004</i>	0.014	0.011
		0.016	0.021	0.020	0.017	0.021	0.021	<i>0.015</i>	0.011
		0.005	0.006	0.004	0.005	0.005	<i>0.004</i>	0.014	0.012
		0.022	0.017	0.027	0.026	<i>0.015</i>	0.030	0.015	0.012
	5	0.031	0.018	0.018	0.019	0.008	<i>0.009</i>	0.058	0.111
		0.060	0.044	0.042	0.091	<i>0.024</i>	0.022	0.060	0.112
		0.035	0.014	0.019	0.039	0.009	<i>0.009</i>	0.056	0.107
		0.062	0.030	0.040	0.079	<i>0.024</i>	0.020	0.064	0.116

Table 3: Mean squared error of $\hat{\beta}_1$ (**minimum**, *next best*).

function	location	lars	enet	glmnet	u.lars	u.enet	u.glmnet	oracular	gwr
step	1	-0.027	-0.024	<i>-0.022</i>	-0.001	0.024	0.025	0.237	-0.333
		-0.080	-0.166	-0.164	0.006	-0.073	<i>-0.073</i>	0.231	-0.329
		-0.029	-0.070	-0.038	0.037	<i>0.030</i>	0.032	0.225	-0.335
		-0.055	-0.163	-0.146	<i>0.027</i>	-0.042	0.007	0.231	-0.330
	2	<i>-0.133</i>	-0.175	-0.174	-0.122	-0.166	-0.167	-0.302	-0.407
		-0.179	-0.201	-0.190	<i>-0.149</i>	-0.153	-0.149	-0.301	-0.404
		<i>-0.125</i>	-0.185	-0.169	-0.110	-0.155	-0.154	-0.305	-0.409
		-0.212	-0.246	-0.248	-0.175	<i>-0.185</i>	-0.202	-0.297	-0.402
	3	-0.003	<i>0.010</i>	0.012	0.019	0.029	0.029	0.060	0.081
		-0.006	-0.080	-0.075	0.039	-0.036	<i>-0.030</i>	0.063	0.084
		-0.021	-0.025	0.004	<i>0.006</i>	0.021	0.024	0.062	0.080
		<i>-0.056</i>	-0.130	-0.110	-0.031	-0.098	-0.092	0.064	0.086
	4	0.074	0.035	<i>0.039</i>	0.093	0.042	0.045	0.300	0.406
		0.079	<i>0.055</i>	0.055	0.094	0.066	0.067	0.303	0.408
		0.046	0.017	<i>0.027</i>	0.062	0.030	0.031	0.307	0.406
		0.087	<i>0.054</i>	0.054	0.101	0.076	0.061	0.303	0.411
	5	-0.005	0.000	0.000	-0.002	0.000	0.000	0.000	0.331
		-0.022	<i>-0.012</i>	-0.013	-0.048	-0.044	-0.043	0.000	0.331
		-0.012	0.000	0.000	-0.053	0.000	0.000	0.000	0.332
		-0.017	-0.011	-0.015	-0.062	<i>-0.006</i>	-0.033	0.000	0.337
gradient	1	-0.320	-0.282	-0.281	0.009	0.009	<i>0.009</i>	0.007	-0.383
		-0.295	-0.286	-0.288	-0.003	0.040	0.039	<i>0.008</i>	-0.380
		-0.312	-0.301	-0.302	-0.001	<i>0.001</i>	0.001	0.003	-0.381
		-0.302	-0.310	-0.312	-0.017	<i>-0.014</i>	-0.049	0.000	-0.388
	2	-0.056	-0.047	-0.046	0.006	0.004	<i>0.004</i>	0.004	-0.068
		-0.046	-0.045	-0.045	0.010	<i>0.009</i>	0.009	0.006	-0.065
		-0.051	-0.053	-0.052	0.003	0.000	<i>0.000</i>	0.001	-0.067
		-0.072	-0.086	-0.089	<i>-0.013</i>	-0.025	-0.038	-0.002	-0.070
	3	0.007	0.004	0.005	0.001	<i>0.002</i>	0.002	0.004	0.015
		0.005	-0.002	-0.002	0.006	<i>0.001</i>	0.001	0.003	0.018
		0.008	0.005	0.007	-0.000	0.001	0.001	<i>-0.000</i>	0.016
		-0.015	-0.048	-0.036	<i>-0.011</i>	-0.039	-0.028	-0.002	0.014
	4	0.053	0.041	0.041	0.002	<i>0.002</i>	0.002	0.003	0.071
		0.035	0.017	0.018	<i>0.001</i>	-0.011	-0.011	-0.001	0.073
		0.051	0.041	0.046	-0.001	<i>-0.000</i>	-0.000	-0.001	0.071
		0.015	-0.022	-0.002	-0.017	-0.028	-0.025	<i>-0.002</i>	0.070
	5	0.315	0.234	0.234	0.007	-0.001	<i>-0.001</i>	0.000	0.388
		0.250	0.203	0.212	0.014	<i>0.005</i>	0.006	0.000	0.387
		0.292	0.239	0.257	0.007	0.002	<i>0.000</i>	0.000	0.384
		0.241	0.173	0.199	<i>-0.009</i>	0.011	0.022	0.000	0.387
parabola	1	0.135	0.049	0.049	0.057	<i>0.017</i>	0.016	0.238	0.334
		0.208	0.166	0.165	0.145	<i>0.107</i>	0.107	0.239	0.335
		0.131	0.039	0.058	0.060	0.013	<i>0.015</i>	0.228	0.326
		0.180	0.112	0.125	0.137	<i>0.064</i>	0.063	0.248	0.334
	2	<i>-0.040</i>	-0.046	-0.046	-0.038	-0.048	-0.048	-0.100	-0.094
		-0.065	-0.072	-0.073	-0.056	<i>-0.062</i>	-0.064	-0.099	-0.093
		<i>-0.046</i>	-0.064	-0.053	-0.041	-0.057	-0.055	-0.106	-0.101
		-0.087	-0.086	-0.078	-0.079	-0.063	<i>-0.063</i>	-0.096	-0.093
	3	<i>-0.052</i>	-0.061	-0.062	-0.045	-0.056	-0.056	-0.162	-0.164
		-0.107	-0.116	-0.114	-0.085	-0.094	<i>-0.093</i>	-0.162	-0.163
		<i>-0.058</i>	-0.077	-0.066	-0.048	-0.064	-0.061	-0.167	-0.171
		-0.121	-0.114	-0.091	-0.099	<i>-0.080</i>	-0.056	-0.161	-0.162
	4	<i>-0.040</i>	-0.051	-0.050	-0.037	-0.051	-0.051	-0.114	-0.102
		<i>-0.082</i>	-0.093	-0.093	-0.078	-0.090	-0.089	-0.115	-0.101
		<i>-0.048</i>	-0.067	-0.057	-0.044	-0.059	-0.058	-0.118	-0.109
		-0.064	-0.079	-0.071	-0.051	<i>-0.060</i>	-0.062	-0.114	-0.100
	5	0.137	0.046	0.044	0.038	<i>0.007</i>	0.007	0.234	0.333
		0.169	0.141	0.129	0.047	0.054	<i>0.050</i>	0.224	0.333
		0.129	0.035	0.049	0.049	0.018	<i>0.021</i>	0.231	0.327
		0.173	0.099	0.121	0.107	0.031	<i>0.053</i>	0.235	0.336

Table 4: Bias of $\hat{\beta}_1$ (**minimum**, *next best*).

function	location	lars	enet	glmnet	u.lars	u.enet	u.glmnet	oracular	gwr
step	1	0.016	0.008	0.008	0.077	0.055	0.055	0.003	0.001
		0.042	0.124	0.125	0.270	0.566	0.585	0.010	0.002
		0.025	0.016	0.011	0.134	0.072	0.068	0.003	0.001
		0.076	0.157	0.151	0.296	0.437	0.444	0.011	0.003
	2	0.008	0.005	0.005	0.006	0.004	0.004	0.001	0.001
		0.025	0.045	0.046	0.022	0.046	0.047	0.003	0.002
		0.008	0.006	0.005	0.007	0.005	0.005	0.001	0.001
		0.020	0.034	0.040	0.018	0.037	0.041	0.003	0.002
	3	0.007	0.005	0.005	0.004	0.002	0.002	0.000	0.001
		0.021	0.050	0.050	0.021	0.063	0.063	0.001	0.002
		0.009	0.008	0.006	0.007	0.006	0.004	0.000	0.001
		0.030	0.047	0.054	0.032	0.054	0.061	0.002	0.002
	4	0.010	0.018	0.018	0.013	0.020	0.020	0.001	0.001
		0.029	0.030	0.030	0.032	0.034	0.033	0.002	0.002
		0.012	0.015	0.018	0.014	0.019	0.019	0.000	0.001
		0.028	0.025	0.029	0.032	0.033	0.031	0.002	0.002
	5	0.005	0.000	0.000	0.031	0.000	0.000	0.000	0.001
		0.021	0.012	0.012	0.103	0.062	0.060	0.000	0.002
		0.007	0.000	0.000	0.046	0.000	0.000	0.000	0.001
		0.020	0.008	0.012	0.127	0.030	0.039	0.000	0.003
gradient	1	0.011	0.011	0.011	0.007	0.006	0.006	0.003	0.001
		0.016	0.018	0.018	0.045	0.082	0.087	0.010	0.002
		0.011	0.008	0.008	0.013	0.004	0.004	0.002	0.001
		0.026	0.043	0.046	0.080	0.167	0.135	0.009	0.002
	2	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.000
		0.003	0.003	0.003	0.003	0.003	0.003	0.002	0.001
		0.002	0.001	0.001	0.001	0.001	0.001	0.001	0.001
		0.008	0.015	0.020	0.008	0.018	0.022	0.002	0.002
	3	0.001	0.001	0.001	0.001	0.000	0.000	0.000	0.000
		0.003	0.009	0.009	0.002	0.011	0.011	0.001	0.001
		0.001	0.001	0.000	0.001	0.000	0.000	0.000	0.001
		0.012	0.022	0.020	0.011	0.023	0.020	0.001	0.002
	4	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.000
		0.006	0.010	0.010	0.003	0.007	0.007	0.002	0.001
		0.002	0.001	0.001	0.002	0.001	0.001	0.001	0.001
		0.013	0.017	0.020	0.010	0.017	0.018	0.002	0.002
	5	0.013	0.025	0.025	0.005	0.003	0.003	0.000	0.001
		0.029	0.029	0.029	0.018	0.012	0.012	0.000	0.002
		0.016	0.022	0.022	0.006	0.004	0.004	0.000	0.001
		0.037	0.026	0.029	0.068	0.012	0.014	0.000	0.003
parabola	1	0.014	0.017	0.017	0.023	0.008	0.008	0.002	0.000
		0.034	0.029	0.030	0.122	0.035	0.038	0.010	0.001
		0.015	0.014	0.017	0.029	0.008	0.008	0.003	0.001
		0.047	0.023	0.026	0.163	0.036	0.039	0.013	0.003
	2	0.002	0.001	0.001	0.003	0.001	0.001	0.000	0.000
		0.010	0.013	0.013	0.012	0.019	0.018	0.002	0.001
		0.002	0.001	0.001	0.003	0.001	0.001	0.000	0.000
		0.015	0.010	0.016	0.018	0.011	0.021	0.002	0.002
	3	0.003	0.001	0.001	0.003	0.001	0.001	0.000	0.000
		0.010	0.011	0.011	0.014	0.015	0.014	0.001	0.001
		0.003	0.001	0.001	0.004	0.001	0.001	0.000	0.000
		0.015	0.008	0.010	0.018	0.009	0.019	0.001	0.002
	4	0.003	0.001	0.001	0.003	0.001	0.001	0.001	0.000
		0.010	0.012	0.012	0.011	0.013	0.013	0.002	0.001
		0.003	0.001	0.001	0.003	0.001	0.001	0.001	0.000
		0.018	0.011	0.022	0.024	0.011	0.026	0.002	0.002
	5	0.013	0.016	0.016	0.018	0.008	0.009	0.003	0.001
		0.032	0.024	0.025	0.090	0.021	0.019	0.010	0.001
		0.018	0.013	0.017	0.037	0.008	0.009	0.003	0.001
		0.032	0.020	0.025	0.068	0.023	0.018	0.009	0.003

Table 5: Variance of $\hat{\beta}_1$ (**minimum**, next best).

function	location	lars	enet	glmnet	u.lars	u.enet	u.glmnet	oracular	gwr
step	1	0.130	0.145	0.145	0.097	0.118	<i>0.118</i>	0.194	0.285
		0.732	0.712	0.664	0.493	<i>0.444</i>	0.415	0.979	1.061
		0.232	0.296	0.302	0.122	0.227	<i>0.225</i>	0.430	0.436
		0.752	0.718	0.710	0.565	0.496	<i>0.515</i>	1.048	1.042
	2	<i>0.222</i>	0.257	0.257	0.213	0.248	0.248	0.319	0.377
		1.082	1.048	1.055	1.070	1.044	<i>1.047</i>	1.221	1.236
		<i>0.304</i>	0.359	0.347	0.291	0.333	0.333	0.437	0.505
		0.984	0.951	0.939	0.954	<i>0.914</i>	0.911	1.087	1.187
	3	<i>0.254</i>	0.264	0.264	0.248	0.260	0.260	0.276	0.269
		0.968	1.006	0.992	0.950	<i>0.924</i>	0.908	1.048	1.057
		<i>0.254</i>	0.284	0.283	0.244	0.276	0.275	0.295	0.291
		<i>0.624</i>	0.701	0.674	0.611	0.654	0.630	0.699	0.693
	4	0.244	0.291	0.291	<i>0.255</i>	0.286	0.286	0.366	0.455
		0.998	1.029	1.047	0.978	<i>0.994</i>	1.018	1.113	1.147
		0.272	0.286	0.285	0.285	0.286	<i>0.281</i>	0.407	0.496
		<i>0.749</i>	0.754	0.751	0.749	0.757	0.751	0.872	0.930
	5	<i>0.223</i>	0.328	0.328	0.162	0.311	0.310	0.337	0.518
		0.651	0.666	0.670	0.505	<i>0.577</i>	0.599	0.789	0.838
		<i>0.204</i>	0.296	0.297	0.148	0.287	0.289	0.296	0.482
		0.930	0.952	0.926	<i>0.817</i>	0.842	0.791	1.166	1.182
gradient	1	0.328	0.300	0.300	0.242	0.236	<i>0.236</i>	0.251	0.375
		0.952	0.847	0.853	0.686	0.664	<i>0.664</i>	0.765	0.939
		0.357	0.359	0.354	0.261	<i>0.266</i>	0.270	0.269	0.372
		0.752	0.786	0.733	0.597	0.750	<i>0.671</i>	0.873	0.842
	2	0.228	0.225	0.225	0.223	0.223	<i>0.223</i>	0.224	0.224
		0.873	0.852	0.852	0.862	0.837	<i>0.837</i>	0.862	0.865
		0.266	0.264	0.264	0.257	0.257	<i>0.257</i>	0.257	0.270
		0.841	0.789	0.798	0.812	0.778	<i>0.781</i>	0.817	0.809
	3	0.248	0.247	0.247	<i>0.246</i>	0.248	0.248	0.249	0.242
		1.140	1.155	1.154	1.141	1.148	1.148	1.147	<i>1.140</i>
		0.317	0.321	0.321	<i>0.313</i>	0.318	0.318	0.318	0.311
		1.249	1.249	<i>1.216</i>	1.248	1.226	1.189	1.297	1.274
	4	0.283	0.291	0.291	0.278	0.281	<i>0.281</i>	0.285	0.299
		0.870	0.871	0.871	0.882	0.866	0.865	0.890	<i>0.866</i>
		0.213	0.209	0.209	0.207	0.208	0.208	<i>0.208</i>	0.214
		0.794	0.818	0.804	0.793	<i>0.793</i>	0.779	0.815	0.820
	5	0.332	0.320	0.320	0.205	0.209	0.209	<i>0.209</i>	0.373
		1.363	1.370	1.353	<i>1.222</i>	1.234	1.195	1.313	1.466
		0.279	0.293	0.298	0.218	0.228	0.228	<i>0.226</i>	0.333
		0.978	1.006	0.992	0.846	0.899	<i>0.889</i>	1.041	1.103
parabola	1	0.151	0.154	0.154	0.126	0.144	<i>0.144</i>	0.212	0.260
		1.192	1.252	1.251	1.078	<i>1.098</i>	1.102	1.414	1.516
		0.286	0.285	0.288	0.229	0.271	<i>0.270</i>	0.367	0.448
		0.848	0.905	0.902	0.753	0.816	<i>0.808</i>	0.968	1.017
	2	<i>0.199</i>	0.201	0.201	0.195	0.201	0.202	0.203	0.201
		1.197	1.176	1.170	1.195	<i>1.161</i>	1.155	1.229	1.223
		0.241	0.252	0.249	<i>0.241</i>	0.246	0.245	0.263	0.263
		1.167	1.129	<i>1.101</i>	1.160	1.118	1.090	1.211	1.207
	3	<i>0.213</i>	0.224	0.223	0.210	0.225	0.224	0.239	0.238
		<i>1.042</i>	1.047	1.048	1.038	1.046	1.044	1.072	1.081
		<i>0.228</i>	0.239	0.238	0.226	0.237	0.237	0.262	0.258
		0.925	0.943	0.938	<i>0.911</i>	0.910	0.920	1.026	1.018
	4	0.248	<i>0.243</i>	0.243	0.248	0.245	0.246	0.246	0.245
		<i>1.354</i>	1.379	1.379	1.344	1.359	1.361	1.421	1.393
		<i>0.291</i>	0.301	0.296	0.291	0.295	0.294	0.317	0.309
		1.019	1.020	1.001	0.995	<i>0.992</i>	0.960	1.086	1.058
	5	0.239	0.236	0.235	0.201	<i>0.235</i>	0.235	0.284	0.342
		<i>0.953</i>	1.007	1.005	0.902	0.962	0.965	1.009	1.059
		<i>0.193</i>	0.229	0.223	0.163	0.219	0.216	0.253	0.295
		<i>0.850</i>	0.905	0.917	0.746	0.891	0.873	0.987	1.044

Table 6: Mean squared error of \hat{Y} (**minimum**, *next best*).