

# Local Adaptive Grouped Regularization and its Oracle Properties

Wesley Brooks, Jun Zhu, Zudi Lu

July 8, 2014

## 1 Introduction

Whereas the coefficients in traditional linear regression are scalar constants, the coefficients in a varying coefficients regression (VCR) model are functions - often *smooth* functions - of some effect-modifying variable ??.

Current practice for VCR models relies on global model selection to decide which variables should be included in the model, meaning that predictors are identified as relevant or irrelevant over the entire domain  $\mathcal{D}$ . ? describe a method for globally selecting the relevant predictors in a VCR model where the coefficient functions are estimated with P-splines. ? show a method for doing global variable selection in a VCR model where the coefficient functions are estimated by basis expansion.

Local adaptive grouped regularization (LAGR) is developed here as a method to select only the locally relevant predictors at any specific location  $\mathbf{s}$  in the

domain  $\mathcal{D}$  of a VCR model. The method of LAGR applies to VCR models where the coefficients are estimated using locally linear kernel smoothing.

Using kernel smoothing for nonparametric regression is described in detail in ?. The extension to estimating VCR models is made by ? for a VCR a univariate effect-modifying variable, and by ? for two-dimensional effect-modifying variable and autocorrelation among the observed response. These methods minimize the boundary effect ? by estimating the coefficients as local polynomials of odd degree (usually locally linear).

For linear regression models, the adaptive lasso (AL) ? produces consistent estimates of the coefficients and has been shown to have appealing properties for automating variable selection, which under suitable conditions include the “oracle” property of asymptotically including exactly the correct set of covariates and estimating their coefficients as well as if the correct covariates were known in advance. For data where the observed variables fall into mutually exclusive groups that are known in advance, the adaptive group lasso ?? has similar oracle properties to the adaptive lasso while doing selection at the level of groups rather than individual variables. The proposed LAGR method uses the adaptive group lasso for local variable selection and coefficient estimation in a locally linear regression model. We show that LAGR possesses the oracle properties of asymptotically selecting exactly the correct local covariates and estimating their local coefficients as accurately as would be possible if the identity of the nonzero coefficients for the local model were known in advance.

The remainder of this document is organized as follows. The kernel-based VCR model is described in Section 2; the proposed LAGR technique and its oracle properties are presented in Section 3; in Section 4, the performance of the proposed LAGR technique is evaluated in a simulation study, and in Section 5

the proposed method is applied to the Boston house price dataset. Proofs of the theorems appear in Appendix A.

## 2 Varying coefficients regression

### 2.1 Model

Consider  $n$  data points, observed at sampling locations  $\mathbf{s}_i = (s_{i,1}, s_{i,2})^T$  for  $i = 1, \dots, n$ , which are distributed in a spatial domain  $\mathcal{D} \subset \mathbb{R}^2$  according to a density  $f(\mathbf{s})$  with  $\mathbf{s} \in \mathcal{D}$ . For  $i = 1, \dots, n$ , let  $y(\mathbf{s}_i)$  and  $\mathbf{x}(\mathbf{s}_i)$  denote, respectively, the univariate response and the  $(p+1)$ -variate vector of covariates measured at location  $\mathbf{s}_i$ . At each location  $\mathbf{s}_i$ , assume that the outcome is related to the covariates by a linear regression where the coefficients  $\boldsymbol{\beta}(\mathbf{s}_i)$  are functions in two dimensions and  $\varepsilon(\mathbf{s}_i)$  is random error at location  $\mathbf{s}_i$ . That is,

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

Further assume that the error term  $\varepsilon(\mathbf{s}_i)$  is normally distributed with zero mean and variance  $\sigma^2$ , and that  $\varepsilon(\mathbf{s}_i)$ ,  $i = 1, \dots, n$  are independent. That is,

$$\boldsymbol{\varepsilon} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2). \quad (2)$$

In the context of nonparametric regression, the boundary-effect bias can be reduced by local polynomial modeling, usually in the form of a locally linear model ?. Here, to prepare for the estimation of locally linear coefficients, we augment the local design matrix with covariate-by-location interactions in two

dimensions ?. The augmented local design matrix at location  $\mathbf{s}_i$  is

$$\mathbf{Z}(\mathbf{s}_i) = (\mathbf{X} \quad \mathbf{L}_i \mathbf{X} \quad \mathbf{M}_i \mathbf{X}) \quad (3)$$

where  $\mathbf{X}$  is the unaugmented matrix of covariates,  $\mathbf{L}_i = \text{diag}\{s_{i',1} - s_{i,1}\}$  and  $\mathbf{M}_i = \text{diag}\{s_{i',2} - s_{i,2}\}$  for  $i' = 1, \dots, n$ .

Now we have that  $Y(\mathbf{s}_i) = \{\mathbf{Z}(\mathbf{s}_i)\}_i^T \boldsymbol{\zeta}(\mathbf{s}_i) + \varepsilon(\mathbf{s}_i)$ , where  $\{\mathbf{Z}(\mathbf{s}_i)\}_i^T$  is the  $i$ th row of the matrix  $\mathbf{Z}(\mathbf{s}_i)$  as a row vector, and  $\boldsymbol{\zeta}(\mathbf{s}_i)$  is the vector of local coefficients at location  $\mathbf{s}_i$ , augmented with the local gradients of the coefficient surfaces in the two spatial dimensions, indicated by  $\nabla_u$  and  $\nabla_v$ :

$$\boldsymbol{\zeta}(\mathbf{s}_i) = (\boldsymbol{\beta}(\mathbf{s}_i)^T, \nabla_u \boldsymbol{\beta}(\mathbf{s}_i)^T, \nabla_v \boldsymbol{\beta}(\mathbf{s}_i)^T)^T$$

## 2.2 Local Likelihood and Coefficient Estimation

The total log-likelihood of the observed data is the sum of the log-likelihood of each individual observation:

$$\ell(\boldsymbol{\zeta}) = -(1/2) \sum_{i=1}^n \left[ \log \sigma^2 + \sigma^{-2} \{y(\mathbf{s}_i) - \mathbf{z}'(\mathbf{s}_i) \boldsymbol{\zeta}(\mathbf{s}_i)\}^2 \right]. \quad (4)$$

Since there are a total of  $n \times 3(p+1) + 1$  parameters for  $n$  observations, the model is not identifiable and it is not possible to directly maximize the total likelihood. But when the coefficient functions are smooth, the coefficients at location  $\mathbf{s}$  can approximate the coefficients within some neighborhood of  $\mathbf{s}$ , with the quality of the approximation declining as the distance from  $\mathbf{s}$  increases.

This intuition is formalized by the local (log-)likelihood, which is maximized at location  $\mathbf{s}$  to estimate the local coefficients  $\boldsymbol{\zeta}(\mathbf{s})$ :

$$\ell\{\boldsymbol{\zeta}(\mathbf{s})\} = -(1/2) \sum_{i=1}^n K_h(\|\mathbf{s} - \mathbf{s}_i\|) \left[ \log \sigma^2 + \sigma^{-2} \{y(\mathbf{s}_i) - \mathbf{z}'(\mathbf{s}_i)\boldsymbol{\zeta}(\mathbf{s})\}^2 \right] \quad (5)$$

where  $h$  is a bandwidth parameter and the  $K_h(\|\mathbf{s} - \mathbf{s}_i\|)$  for  $i = 1, \dots, n$  are local weights from a kernel function. For instance, the Epanechnikov kernel is defined as ?:

$$K_h(\|\mathbf{s}_i - \mathbf{s}_{i'}\|) = h^{-2} K(h^{-1}\|\mathbf{s}_i - \mathbf{s}_{i'}\|)$$

$$K(x) = \begin{cases} (3/4)(1 - x^2) & \text{if } x < 1, \\ 0 & \text{if } x \geq 1. \end{cases} \quad (6)$$

Letting  $\mathbf{W}(\mathbf{s}) = \text{diag}\{K_h(\|\mathbf{s} - \mathbf{s}_i\|)\}$  be a diagonal matrix of kernel weights, the local likelihood is maximized by weighted least squares:

$$\mathcal{S}\{\boldsymbol{\zeta}(\mathbf{s})\} = (1/2) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\}^T \mathbf{W}(\mathbf{s}) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\}^T$$

Thus, we have

$$\tilde{\boldsymbol{\zeta}}(\mathbf{s}) = \{\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}(\mathbf{s})\}^{-1} \mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Y}$$

Now Theorem 3 of ? says that, for any given  $\mathbf{s}$

$$\sqrt{nh^2 f(\mathbf{s})} \left[ \hat{\beta}(\mathbf{s}) - \beta(\mathbf{s}) - (1/2)\kappa_0^{-1}\kappa_2 h^2 \{\beta_{uu}(\mathbf{s}) + \beta_{vv}(\mathbf{s})\} \right] \xrightarrow{D} N(\mathbf{0}, \kappa_0^{-2}\nu_0\sigma^2\Psi^{-1})$$

### 3 Local Variable Selection with LAGR

#### 3.1 The LAGR-Penalized Local Likelihood

Estimating the local coefficients by (??) relies on *a priori* variable selection. A new method of penalized regression to simultaneously select the locally relevant predictors and estimate the local coefficients. For this purpose, each raw covariate is grouped with its covariate-by-location interactions. That is,  $\zeta_j(\mathbf{s}) = (\beta_j(\mathbf{s}) \ \nabla_u \beta_j(\mathbf{s}) \ \nabla_v \beta_j(\mathbf{s}))^T$  for  $j = 1, \dots, p$ . By the mechanism of the group lasso, variables within the same group are included in or dropped from the model together. The intercept group is left unpenalized. The proposed LAGR penalty is an adaptive  $\ell_1$  penalty akin to the adaptive group lasso ??.

More specifically, we consider the penalized local sum of squares at location  $\mathbf{s}$ :

$$\mathcal{J}\{\zeta(\mathbf{s})\} = \mathcal{S}\{\zeta(\mathbf{s})\} + \mathcal{P}\{\zeta(\mathbf{s})\}$$

where  $\mathcal{S}\{\zeta(\mathbf{s})\} = (1/2) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\zeta(\mathbf{s})\}^T \mathbf{W}(\mathbf{s}) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\zeta(\mathbf{s})\}^T$  is the locally weighted sum of squares,  $\mathcal{P}\{\zeta(\mathbf{s})\} = \sum_{j=1}^p \phi_j(\mathbf{s}) \|\zeta_j(\mathbf{s})\|$  is a local adaptive grouped regularization (LAGR) penalty, and  $\|\cdot\|$  is the  $L_2$ -norm.

The LAGR penalty for the  $j$ th group of coefficients  $\zeta_j(\mathbf{s})$  at location  $\mathbf{s}$  is  $\phi_j(\mathbf{s}) = \lambda_n(\mathbf{s}) \|\tilde{\zeta}_j(\mathbf{s})\|^{-\gamma}$ , where  $\lambda_n(\mathbf{s}) > 0$  is a local tuning parameter applied to all

coefficients at location  $\mathbf{s}$  and  $\tilde{\boldsymbol{\zeta}}_j(\mathbf{s})$  is the vector of unpenalized local coefficients from (??).

### 3.2 Oracle properties of LAGR

**Theorem 1** (Asymptotic normality). *If  $h^{-1}n^{-1/2}a_n \xrightarrow{p} 0$  and  $hn^{-1/2}b_n \xrightarrow{p} \infty$  then*

$$h\sqrt{n} \left[ \hat{\boldsymbol{\beta}}_{(a)}(\mathbf{s}) - \boldsymbol{\beta}_{(a)}(\mathbf{s}) - \frac{\kappa_2 h^2}{2\kappa_0} \{ \nabla_{uu}^2 \boldsymbol{\beta}_{(a)}(\mathbf{s}) + \nabla_{vv}^2 \boldsymbol{\beta}_{(a)}(\mathbf{s}) \} \right] \xrightarrow{d} N(0, f(\mathbf{s})^{-1} \kappa_0^{-2} \nu_0 \sigma^2 \Psi^{-1})$$

**Theorem 2** (Selection consistency). *If  $h^{-1}n^{-1/2}a_n \xrightarrow{p} \infty$  and  $hn^{-1/2}b_n \xrightarrow{p} \infty$  then  $P \left\{ \|\hat{\boldsymbol{\zeta}}_j(\mathbf{s})\| = 0 \right\} \rightarrow 0$  if  $j \leq p_0$  and  $P \left\{ \|\hat{\boldsymbol{\zeta}}_j(\mathbf{s})\| = 0 \right\} \rightarrow 1$  if  $j > p_0$ .*

**Remarks** Together, Theorem 1 and Theorem 2 indicate that the LAGR estimates have the same asymptotic distribution as a local regression model where the nonzero coefficients are known in advance ?, and that the LAGR estimates of true zero coefficients go to zero with probability one. Thus, selection and estimation by LAGR has the oracle property.

**A note on rates** To prove the oracle properties of LAGR, we assumed that  $h^{-1}n^{-1/2}a_n \xrightarrow{p} 0$  and  $hn^{-1/2}b_n \xrightarrow{p} \infty$ . Therefore,  $h^{-1}n^{-1/2}\lambda_n(\mathbf{s}) \rightarrow 0$  for  $j \leq p_0$  and  $hn^{-1/2}\lambda_n(\mathbf{s})\|\boldsymbol{\zeta}_j(\mathbf{s})\|^{-\gamma} \rightarrow \infty$  for  $j > p_0$ .

We require that  $\lambda_n(\mathbf{s})$  can satisfy both assumptions. Suppose  $\lambda_n(\mathbf{s}) = n^\alpha$ , and recall that  $h = O(n^{-1/6})$  and  $\|\tilde{\boldsymbol{\zeta}}_p(\mathbf{s})\| = O(h^{-1}n^{-1/2})$ . Then  $h^{-1}n^{-1/2}\lambda_n(\mathbf{s}) = O(n^{-1/3+\alpha})$  and  $hn^{-1/2}\lambda_n(\mathbf{s})\|\tilde{\boldsymbol{\zeta}}_p(\mathbf{s})\|^{-\gamma} = O(n^{-2/3+\alpha+\gamma/3})$ .

So  $(2 - \gamma)/3 < \alpha < 1/3$ , which can only be satisfied for  $\gamma > 1$ .

### 3.3 Selecting the tuning parameter $\lambda_n(\mathbf{s})$

In practical application, it is necessary to select the LAGR tuning parameter  $\lambda_n(\mathbf{s})$  for each local model. A popular approach in other lasso-type problems is to select the tuning parameter that maximizes a criterion that approximates the expected log-likelihood of a new, independent data set drawn from the same distribution. This is the framework of Mallows' Cp [1], Stein's unbiased risk estimate (SURE) [2] and Akaike's information criterion (AIC) [3].

These criteria use a so-called covariance penalty [4] to estimate the bias due to using the same data set to select a model and to estimate its parameters. We adopt the approximate degrees of freedom for the adaptive group lasso from [5] and minimize the AICc [6] to select the tuning parameter  $\lambda_n(\mathbf{s})$ :

$$\begin{aligned} \hat{df}(\lambda; \mathbf{s}) &= \sum_{j=1}^p I\left(\|\hat{\boldsymbol{\zeta}}(\lambda; \mathbf{s})\| > 0\right) + \sum_{j=1}^p \frac{\|\hat{\boldsymbol{\zeta}}(\lambda; \mathbf{s})\|}{\|\tilde{\boldsymbol{\zeta}}(\mathbf{s})\|} (p_j - 1) \\ \text{AIC}_c(\lambda; \mathbf{s}) &= \sum_{i=1}^n K_h(\|\mathbf{s} - \mathbf{s}_i\|) \sigma^{-2} \left\{ y(\mathbf{s}_i) - z'(\mathbf{s}_i) \hat{\boldsymbol{\zeta}}(\lambda; \mathbf{s}) \right\}^2 + 2\hat{df}(\lambda; \mathbf{s}) + \frac{2\hat{df}(\lambda; \mathbf{s}) \left\{ \hat{df}(\lambda; \mathbf{s}) + 1 \right\}}{\sum_{i=1}^n K_h(\|\mathbf{s} - \mathbf{s}_i\|) - \hat{df}(\lambda; \mathbf{s}) - 1} \end{aligned}$$

where the local coefficient estimate has been written  $\hat{\boldsymbol{\zeta}}(\lambda; \mathbf{s})$  to emphasize that it depends on the tuning parameter.



## 4 Simulations

### 4.1 Simulation Setup

A simulation study was conducted to assess the performance of the method described in Sections 2–3. Data were simulated on the domain  $[0, 1]^2$ , which was divided into a  $30 \times 30$  grid. Each of  $p = 5$  covariates  $X_1, \dots, X_5$  was simulated by a Gaussian random field with mean zero and exponential covariance function  $\text{Cov}(X_{ji}, X_{ji'}) = \sigma_x^2 \exp(-\tau_x^{-1} \delta_{ii'})$  where  $\sigma_x^2 = 1$  is the variance,  $\tau_x = 0.1$  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 matrix  $\mathbf{X} = (X_1 \cdots X_5)$  by  $\mathbf{R}$ , where  $\mathbf{R}$  is the Cholesky decomposition of the covariance matrix  $\mathbf{\Sigma} = \mathbf{R}'\mathbf{R}$ . The covariance matrix  $\mathbf{\Sigma}$  is a  $5 \times 5$  matrix that has ones on the diagonal and  $\rho$  for all off-diagonal entries, where  $\rho$  is the between-covariate correlation.

The simulated response was  $y_i = \mathbf{x}_i' \boldsymbol{\beta}_i + \varepsilon_i$  for  $i = 1, \dots, n$  where  $n = 900$  and the  $\varepsilon_i$ 's were iid Gaussian with mean zero and variance  $\sigma_\varepsilon^2$ . The simulated data included the response  $y$  and five covariates  $x_1, \dots, x_5$ . The true data-generating model uses only  $x_1$ . The variables  $x_2, \dots, x_5$  are included to assess performance in model selection.

Three different functions were used for the coefficient surface  $\beta_1(\mathbf{s})$ . They are plotted in Figure 1, and their mathematical forms are listed in (7). The first is a step function, which is equal to zero in 40% of the spatial domain, equal to one in a different 40% of the spatial domain, and increases linearly in the middle 20% of the domain. The second is a gradient function, which increases linearly from zero at one end of the domain to one at the other. The final coefficient

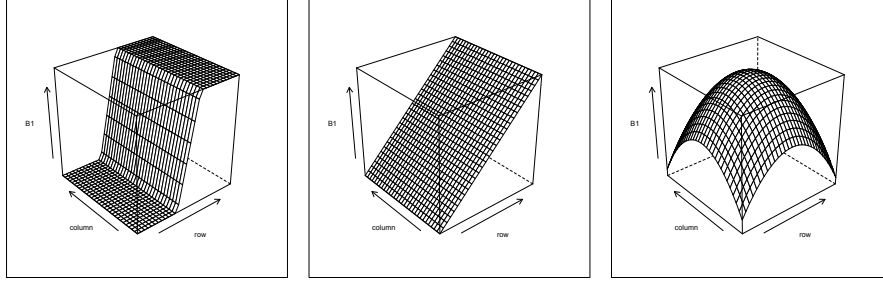


Figure 1: These are, respectively, the step, gradient, and parabola functions that were used for the coefficient function  $\beta_1(\mathbf{s})$  in the VCR model  $y(\mathbf{s}_i) = x_1(\mathbf{s}_i)\beta_1(\mathbf{s}_i) + \varepsilon(\mathbf{s}_i)$  when generating the data for the simulation study.

function is a parabola taking its maximum value of 1 at the center of the domain and falling to zero at each corner of the domain.

$$\beta_{step}(\mathbf{s}) = \begin{cases} 1 & \text{if } s_x > 0.6 \\ 5s_x - 2 & \text{if } 0.4 < s_x \leq 0.6 \\ 0 & \text{o.w.} \end{cases}$$

$$\beta_{gradient}(\mathbf{s}) = s_x$$

$$\beta_{parabola}(\mathbf{s}) = 1 - \frac{(s_x - 0.5)^2 + (s_y - 0.5)^2}{0.5} \quad (7)$$

In total, three parameters were varied to produce 18 settings, each of which was simulated 100 times. There were three functional forms for the coefficient surface  $\beta_1(\mathbf{s})$ ; data was simulated both with low ( $\rho = 0$ ), medium ( $\rho = 0.5$ ), and high ( $\rho = 0.9$ ) 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 error term. The simulation settings are enumerated in Table ??.

Simulation settings			MISE			MISE	
$\beta_1(\mathbf{s})$	$\rho$	$\sigma_\varepsilon$	LAGR	$\hat{\beta}_1$ VCR	Oracle	$\hat{\beta}_2, \dots, \hat{\beta}_5$ LAGR	VCR
step	0	0.5	<i>0.02</i>	0.02	<b>0.01</b>	<b>0.00</b>	0.01
		1.0	<i>0.03</i>	0.03	<b>0.02</b>	<b>0.00</b>	0.02
	0.5	0.5	<i>0.02</i>	0.02	<b>0.01</b>	<b>0.00</b>	0.01
		1.0	<i>0.03</i>	0.05	<b>0.02</b>	<b>0.00</b>	0.03
	0.9	0.5	<i>0.03</i>	0.05	<b>0.01</b>	<b>0.00</b>	0.04
		1.0	<i>0.12</i>	0.17	<b>0.02</b>	<b>0.02</b>	0.15
gradient	0	0.5	0.01	<i>0.01</i>	<b>0.00</b>	<b>0.00</b>	0.00
		1.0	0.03	<i>0.02</i>	<b>0.01</b>	<b>0.00</b>	0.02
	0.5	0.5	0.01	<i>0.01</i>	<b>0.00</b>	<b>0.00</b>	0.01
		1.0	0.04	<i>0.03</i>	<b>0.01</b>	<b>0.00</b>	0.03
	0.9	0.5	<i>0.03</i>	0.04	<b>0.00</b>	<b>0.00</b>	0.04
		1.0	<i>0.14</i>	0.14	<b>0.01</b>	<b>0.02</b>	0.15
parabola	0	0.5	0.01	<i>0.01</i>	<b>0.01</b>	<b>0.00</b>	0.00
		1.0	0.03	<i>0.02</i>	<b>0.02</b>	<b>0.00</b>	0.02
	0.5	0.5	0.01	<i>0.01</i>	<b>0.01</b>	<b>0.00</b>	0.01
		1.0	0.03	<i>0.03</i>	<b>0.02</b>	<b>0.00</b>	0.03
	0.9	0.5	<i>0.02</i>	0.04	<b>0.01</b>	<b>0.00</b>	0.04
		1.0	0.17	<i>0.14</i>	<b>0.02</b>	<b>0.03</b>	0.15

Table 1: Listing of the simulation settings used to assess the performance of LAGR models versus oracle selection and no selection.

## 4.2 Simulation Results

The results are presented in terms of the mean integrated squared error (MISE) of the coefficient surface estimates  $\hat{\beta}_1(\mathbf{s}), \dots, \hat{\beta}_5(\mathbf{s})$ , the MISE of the fitted response  $\hat{y}(\mathbf{s})$ , and the frequency with which the coefficient surface estimates  $\hat{\beta}_1(\mathbf{s}), \dots, \hat{\beta}_5(\mathbf{s})$  estimated by LAGR were zero. The performance of LAGR was compared to that of a VCR model without variable selection, and to a VCR model with oracular selection. Oracular selection means that exactly the correct set of covariates was used to fit each local model.

The MISE of the estimates of  $\beta_1(\mathbf{s})$  are in Table ???. Recall that  $\beta_2(\mathbf{s}), \dots, \beta_5(\mathbf{s})$  are exactly zero across the entire domain. Oracle selection will estimate these coefficients perfectly, so we focus on the comparison between estimation by

Simulation settings			Zero frequency	MISE		
$\beta_1(\mathbf{s})$	$\rho$	$\sigma_\varepsilon$	$\hat{\beta}_2, \dots, \hat{\beta}_5$	LAGR	$\hat{y}$ VCR	Oracle
step	0	0.5	0.97	0.25	0.26	<b>0.25</b>
		1.0	0.96	1.00	<b>1.00</b>	0.99
	0.5	0.5	0.96	0.26	0.26	<b>0.25</b>
		1.0	0.92	0.99	<b>1.00</b>	0.98
	0.9	0.5	0.86	0.27	0.30	<b>0.25</b>
		1.0	0.85	1.08	1.14	<b>0.98</b>
gradient	0	0.5	0.96	0.25	<b>0.25</b>	0.25
		1.0	0.95	<b>0.99</b>	0.99	0.97
	0.5	0.5	0.94	0.25	<b>0.25</b>	0.24
		1.0	0.92	1.00	<b>1.00</b>	0.97
	0.9	0.5	0.80	0.27	0.28	<b>0.24</b>
		1.0	0.85	1.09	1.12	<b>0.97</b>
parabola	0	0.5	0.97	0.25	<b>0.25</b>	0.25
		1.0	0.94	<b>1.00</b>	1.00	0.98
	0.5	0.5	0.95	<b>0.25</b>	0.25	0.25
		1.0	0.88	<b>1.00</b>	1.00	0.97
	0.9	0.5	0.79	0.26	0.28	<b>0.24</b>
		1.0	0.78	1.13	1.12	<b>0.98</b>

Table 2: The MISE for the fitted output in each simulation setting, under variable selection via LAGR, no variable selection, and oracular variable selection. Highlighting indicates the **closest** and *next-closest* to the actual error variance  $\sigma_\varepsilon^2$  for that setting.

LAGR and by the VCR model with no selection. The MISE of the estimates of  $\beta_2(\mathbf{s}), \dots, \beta_5(\mathbf{s})$  for each simulation setting are enumerated in Table ??, which shows that for every simulation setting, LAGR selection and estimation is more accurate than the standard VCR model.

From Table 2 we see that LAGR has good ability to identify zero-coefficient covariates. The frequency with which  $\beta_2(\mathbf{s}), \dots, \beta_5(\mathbf{s})$  were dropped from the LAGR models ranged from 0.78 to 0.97.

The MISE of the fitted  $\hat{y}(\mathbf{s})$  is listed in Table 2, where the highlighting is based on which methods estimate an error variance that is closest to the known truth for the simulation. The results are all very similar to each other, indicating that

no method was consistently better than the others in this simulation at fitting the model output.

### 4.3 Discussion

The proposed LAGR method was accurate in selection and estimation, with estimation accuracy for  $\beta_1(\mathbf{s})$  about equal to that of the VCR model with no selection, and with consistently better accuracy for estimating  $\beta_2(\mathbf{s}), \dots, \beta_5(\mathbf{s})$ .

There was minimal difference in the performance of the proposed LAGR method between low ( $\sigma_\varepsilon = 0.5$ ) and high ( $\sigma_\varepsilon = 1$ ) error variance, and between no ( $\rho = 0$ ) and moderate ( $\rho = 0.5$ ) correlation among the covariates. But the selection and estimation accuracy did decline when there was high ( $\rho = 0.9$ ) correlation among the predictor variables.

## 5 Data example

The proposed LAGR estimation method was used to estimate the coefficients in a VCR model of the effect of some covariates on the price of homes in Boston ????. The data source is based on the 1970 U.S. census. In the data, we have the median price of homes sold in 506 census tracts (MEDV), along with some potential predictor variables. The predictor variables are CRIM (the per-capita crime rate in the tract), RM (the mean number of rooms for houses sold in the tract), RAD (an index of how accessible the tract is from Boston’s radial roads), TAX (the property tax per \$10,000 of property value), and LSTAT (the percentage of the tract’s residents who are considered “lower status”).

The bandwidth parameter was set to 0.2 for a nearest neighbors-type bandwidth, meaning that the sum of kernel weights for each local model was 20% of the total number of observations. The kernel used was the Epanechnikov kernel.

A summary of the local coefficients is in Table 3. It indicates that RM is the only predictor variable with a positive mean of the local coefficients, but also that the mean of the local coefficients of RM is the largest coefficient - at 1.92, it is more than twice as large in magnitude as the mean local coefficient of LSTAT ( $-0.72$ ), which is second-largest. The coefficient of the CRIM variable was estimated to be exactly zero at 49% of the locations. The percentage for the RAD variable was 37%.

Estimates of the regression coefficients are plotted in Figure 2. One interesting result is that LAGR indicates that the TAX variable was nowhere an important predictor of the median house price. Another is that the coefficients of CRIM and LSTAT are everywhere negative or zero (meaning that a greater crime rate or proportion of lower-status individuals is associated with a lower median house price where the effect is discernable) and that of RM is positive (meaning that a greater average number of rooms per house is associated with a greater median house price). The coefficient of RAD is positive in some areas and negative in others. This indicates that there are parts of Boston where improved access to radial roads is associated with a greater median house price and parts where it is associated with a lesser median house price.

In their example using the same data, ? estimated that the coefficients of RAD and LSTAT should be constant, at 0.36 and -0.45, respectively. That conclusion differs from our result, which says that the mean local coefficient of RAD is actually negative ( $-0.08$ ), while our mean fitted local coefficient for LSTAT was more negative than the estimate of ?.

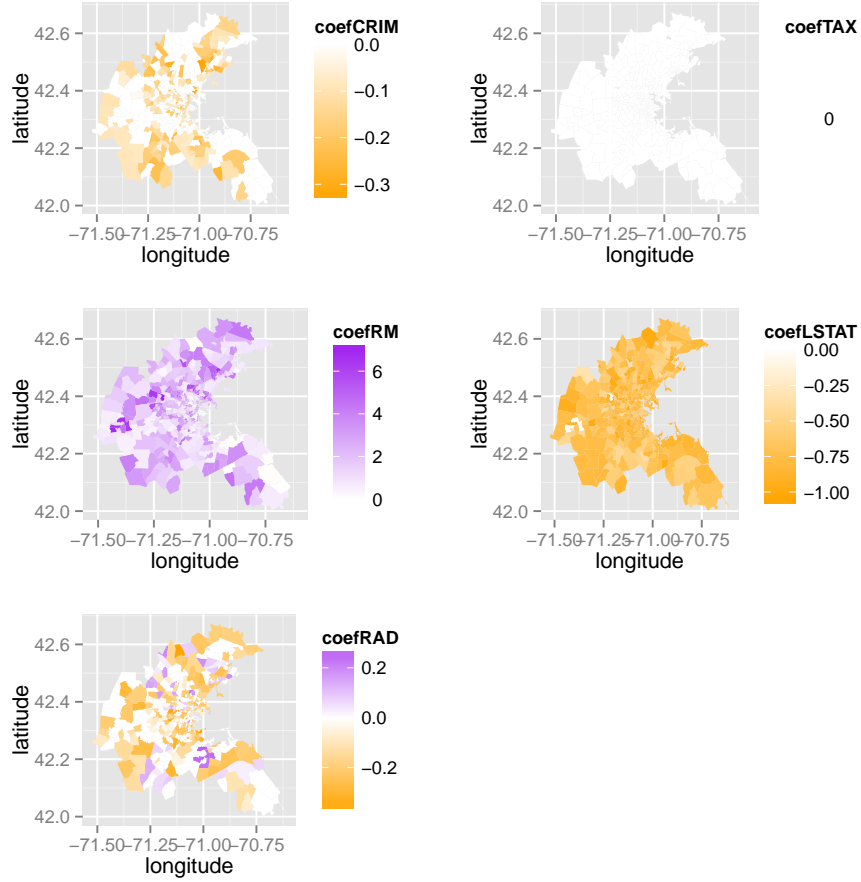


Figure 2: The LAGR estimates of coefficients for the Boston house price data.

	Mean	SD	Prop. zero
CRIM	-0.07	0.08	0.49
RM	1.92	1.43	0.02
RAD	-0.08	0.13	0.37
TAX	0.00	0.00	1.00
LSTAT	-0.72	0.16	0.01

Table 3: The mean, standard deviation, and proportion of zeros among the local coefficients in a model for the median house price in census tracts in Boston, with coefficients selected and fitted by LAGR.

## A Proofs of theorems

*Proof of theorem 1.*

□

Define  $V_4^{(n)}(\mathbf{u})$  to be the

$$\begin{aligned}
V_4^{(n)}(\mathbf{u}) &= \mathcal{J} \left\{ \boldsymbol{\zeta}(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u} \right\} - \mathcal{J} \left\{ \boldsymbol{\zeta}(\mathbf{s}) \right\} \\
&= (1/2) \left[ \mathbf{Y} - \mathbf{Z}(\mathbf{s}) \left\{ \boldsymbol{\zeta}(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u} \right\} \right]^T \mathbf{W}(\mathbf{s}) \left[ \mathbf{Y} - \mathbf{Z}(\mathbf{s}) \left\{ \boldsymbol{\zeta}(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u} \right\} \right] \\
&\quad + \sum_{j=1}^p \phi_j(\mathbf{s}) \|\boldsymbol{\zeta}_j(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u}_j\| \\
&\quad - (1/2) \left\{ \mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s}) \right\}^T \mathbf{W}(\mathbf{s}) \left\{ \mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s}) \right\} - \sum_{j=1}^p \phi_j(\mathbf{s}) \|\boldsymbol{\zeta}_j(\mathbf{s})\| \\
&= (1/2) \mathbf{u}^T \left\{ h^{-2}n^{-1} \mathbf{Z}^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \mathbf{Z}(\mathbf{s}) \right\} \mathbf{u} - \mathbf{u}^T \left[ h^{-1}n^{-1/2} \mathbf{Z}^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \left\{ \mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s}) \right\} \right] \\
&\quad + \sum_{j=1}^p n^{-1/2} \phi_j(\mathbf{s}) n^{1/2} \left\{ \|\boldsymbol{\zeta}_j(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u}_j\| - \|\boldsymbol{\zeta}_j(\mathbf{s})\| \right\} \tag{8}
\end{aligned}$$

Note the different limiting behavior of the third term between the cases  $j \leq p_0$

and  $j > p_0$ :

**Case  $j \leq p_0$**  If  $j \leq p_0$  then  $n^{-1/2}\phi_j(\mathbf{s}) \rightarrow n^{-1/2}\lambda_n(\mathbf{s})\|\boldsymbol{\zeta}_j(\mathbf{s})\|^{-\gamma}$  and  $|\sqrt{n} \{ \|\boldsymbol{\zeta}_j(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u}_j\| - \|\boldsymbol{\zeta}_j(\mathbf{s})\| \}|$   
 $h^{-1}\|\mathbf{u}_j\|$  so

$$\lim_{n \rightarrow \infty} \phi_j(\mathbf{s}) \left( \|\boldsymbol{\zeta}_j(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u}_j\| - \|\boldsymbol{\zeta}_j(\mathbf{s})\| \right) \leq h^{-1}n^{-1/2}\phi_j(\mathbf{s})\|\mathbf{u}_j\| \leq h^{-1}n^{-1/2}a_n\|\mathbf{u}_j\| \rightarrow 0$$

**Case  $j > p_0$**  If  $j > p_0$  then  $\phi_j(\mathbf{s}) \left( \|\boldsymbol{\zeta}_j(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u}_j\| - \|\boldsymbol{\zeta}_j(\mathbf{s})\| \right) = \phi_j(\mathbf{s})h^{-1}n^{-1/2}\|\mathbf{u}_j\|$ .

And note that  $h = O(n^{-1/6})$  so that if  $hn^{-1/2}b_n \xrightarrow{p} \infty$  then  $h^{-1}n^{-1/2}b_n \xrightarrow{p} \infty$ .



Now, if  $\|\mathbf{u}_j\| \neq 0$  then

$$h^{-1}n^{-1/2}\phi_j(\mathbf{s})\|\mathbf{u}_j\| \geq h^{-1}n^{-1/2}b_n\|\mathbf{u}_j\| \rightarrow \infty$$

. On the other hand, if  $\|\mathbf{u}_j\| = 0$  then  $h^{-1}n^{-1/2}\phi_j(\mathbf{s})\|\mathbf{u}_j\| = 0$ .

Thus, the limit of  $V_4^{(n)}(\mathbf{u})$  is the same as the limit of  $V_4^{*(n)}(\mathbf{u})$  where

$$V_4^{*(n)}(\mathbf{u}) = \begin{cases} (1/2)\mathbf{u}^T \{h^{-2}n^{-1}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}(\mathbf{s})\} \mathbf{u} - \mathbf{u}^T [h^{-1}n^{-1/2}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\zeta(\mathbf{s})\}] & \text{if } \|\mathbf{u}_j\| = 0 \ \forall j > t \\ \infty & \text{otherwise} \end{cases}$$

From which it is clear that  $V_4^{*(n)}(\mathbf{u})$  is convex and its unique minimizer is  $\hat{\mathbf{u}}^{(n)}$ :

$$\begin{aligned} 0 &= \{h^{-2}n^{-1}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}(\mathbf{s})\} \hat{\mathbf{u}}^{(n)} - [h^{-1}n^{-1/2}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\zeta(\mathbf{s})\}] \\ \therefore \hat{\mathbf{u}}^{(n)} &= \{n^{-1}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}(\mathbf{s})\}^{-1} [hn^{-1/2}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\zeta(\mathbf{s})\}] \end{aligned} \quad (9)$$

By the epiconvergence results of ? and ?, the minimizer of the limiting function is the limit of the minimizers  $\hat{\mathbf{u}}^{(n)}$ . And since, by Lemma 2 of ?,

$$\hat{\mathbf{u}}^{(n)} \xrightarrow{d} N\left(\frac{\kappa_2 h^2}{2\kappa_0} \{\nabla_{uu}^2 \zeta_j(\mathbf{s}) + \nabla_{vv}^2 \zeta_j(\mathbf{s})\}, f(\mathbf{s})\kappa_0^{-2}\nu_0\sigma^2\Psi^{-1}\right) \quad (10)$$

the result is proven.

*Proof of theorem 2.* We showed in Theorem 1 that  $\hat{\zeta}_j(\mathbf{s}) \xrightarrow{p} \zeta_j(\mathbf{s}) + \frac{\kappa_2 h^2}{2\kappa_0} \{\nabla_{uu}^2 \zeta_j(\mathbf{s}) + \nabla_{vv}^2 \zeta_j(\mathbf{s})\}$ , so to complete the proof of selection consistency, it only remains to show that  $P\{\hat{\zeta}_j(\mathbf{s}) = 0\} \rightarrow 1$  if  $j > p_0$ .  $\square$

The proof is by contradiction. Without loss of generality we consider only the case  $j = p$ .

Assume  $\|\hat{\zeta}_p(\mathbf{s})\| \neq 0$ . Then  $Q\{\zeta(\mathbf{s})\}$  is differentiable w.r.t.  $\zeta_p(\mathbf{s})$  and is minimized where

$$\begin{aligned} 0 &= \mathbf{Z}_p^T(\mathbf{s})\mathbf{W}(\mathbf{s})\left\{\mathbf{Y} - \mathbf{Z}_{-p}(\mathbf{s})\hat{\zeta}_{-p}(\mathbf{s}) - \mathbf{Z}_p(\mathbf{s})\hat{\zeta}_p(\mathbf{s})\right\} - \phi_p(\mathbf{s})\frac{\hat{\zeta}_p(\mathbf{s})}{\|\hat{\zeta}_p(\mathbf{s})\|} \\ &= \mathbf{Z}_p^T(\mathbf{s})\mathbf{W}(\mathbf{s})\left[\mathbf{Y} - \mathbf{Z}(\mathbf{s})\zeta(\mathbf{s}) - \frac{h^2\kappa_2}{2\kappa_0}\{\nabla_{uu}^2\zeta(\mathbf{s}) + \nabla_{vv}^2\zeta(\mathbf{s})\}\right] \\ &\quad + \mathbf{Z}_p^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}_{-p}(\mathbf{s})\left[\zeta_{-p}(\mathbf{s}) + \frac{h^2\kappa_2}{2\kappa_0}\{\nabla_{uu}^2\zeta_{-p}(\mathbf{s}) + \nabla_{vv}^2\zeta_{-p}(\mathbf{s})\} - \hat{\zeta}_{-p}(\mathbf{s})\right] \\ &\quad + \mathbf{Z}_p^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}_p(\mathbf{s})\left[\zeta_p(\mathbf{s}) + \frac{h^2\kappa_2}{2\kappa_0}\{\nabla_{uu}^2\zeta_p(\mathbf{s}) + \nabla_{vv}^2\zeta_p(\mathbf{s})\} - \hat{\zeta}_p(\mathbf{s})\right] \\ &\quad - \phi_p(\mathbf{s})\frac{\hat{\zeta}_p(\mathbf{s})}{\|\hat{\zeta}_p(\mathbf{s})\|} \end{aligned} \tag{11}$$

So

$$\begin{aligned} \frac{h}{\sqrt{n}}\phi_p(\mathbf{s})\frac{\hat{\zeta}_p(\mathbf{s})}{\|\hat{\zeta}_p(\mathbf{s})\|} &= \mathbf{Z}_p^T(\mathbf{s})\mathbf{W}(\mathbf{s})\frac{h}{\sqrt{n}}\left[\mathbf{Y} - \mathbf{Z}(\mathbf{s})\zeta(\mathbf{s}) - \frac{h^2\kappa_2}{2\kappa_0}\{\nabla_{uu}^2\zeta(\mathbf{s}) + \nabla_{vv}^2\zeta(\mathbf{s})\}\right] \\ &\quad + \{n^{-1}\mathbf{Z}_p^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}_{-p}(\mathbf{s})\}h\sqrt{n}\left[\zeta_{-p}(\mathbf{s}) + \frac{h^2\kappa_2}{2\kappa_0}\{\nabla_{uu}^2\zeta_{-p}(\mathbf{s}) + \nabla_{vv}^2\zeta_{-p}(\mathbf{s})\} - \hat{\zeta}_{-p}(\mathbf{s})\right] \\ &\quad + \{n^{-1}\mathbf{Z}_p^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}_p(\mathbf{s})\}h\sqrt{n}\left[\zeta_p(\mathbf{s}) + \frac{h^2\kappa_2}{2\kappa_0}\{\nabla_{uu}^2\zeta_p(\mathbf{s}) + \nabla_{vv}^2\zeta_p(\mathbf{s})\} - \hat{\zeta}_p(\mathbf{s})\right] \end{aligned} \tag{12}$$

From Lemma 2 of ?,  $\{n^{-1}\mathbf{Z}_p^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}_{-p}(\mathbf{s})\} = O_p(1)$  and  $\{n^{-1}\mathbf{Z}_p^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}_p(\mathbf{s})\} =$

$O_p(1)$ .

From Theorem 3 of ?, we have that  $h\sqrt{n} \left[ \hat{\zeta}_{-p}(\mathbf{s}) - \zeta_{-p}(\mathbf{s}) - \frac{h^2\kappa_2}{2\kappa_0} \{ \nabla_{uu}^2 \zeta_{-p}(\mathbf{s}) + \nabla_{vv}^2 \zeta_{-p}(\mathbf{s}) \} \right] = O_p(1)$  and  $h\sqrt{n} \left[ \hat{\zeta}_p(\mathbf{s}) - \zeta_p(\mathbf{s}) - \frac{h^2\kappa_2}{2\kappa_0} \{ \nabla_{uu}^2 \zeta_p(\mathbf{s}) + \nabla_{vv}^2 \zeta_p(\mathbf{s}) \} \right] = O_p(1)$ .

So the second and third terms of the sum in (12) are  $O_p(1)$ .

We showed in the proof of 1 that  $h\sqrt{n} \mathbf{Z}_p^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \left[ \mathbf{Y} - \mathbf{Z}(\mathbf{s}) \boldsymbol{\zeta}(\mathbf{s}) - \frac{h^2\kappa_2}{2\kappa_0} \{ \nabla_{uu}^2 \boldsymbol{\zeta}(\mathbf{s}) + \nabla_{vv}^2 \boldsymbol{\zeta}(\mathbf{s}) \} \right] = O_p(1)$ .

The three terms of the sum to the right of the equals sign in (12) are  $O_p(1)$ , so for  $\hat{\zeta}_p(\mathbf{s})$  to be a solution, we must have that  $hn^{-1/2} \phi_p(\mathbf{s}) \hat{\zeta}_p(\mathbf{s}) / \|\hat{\zeta}_p(\mathbf{s})\| = O_p(1)$ .

But since by assumption  $\hat{\zeta}_p(\mathbf{s}) \neq 0$ , there must be some  $k \in \{1, \dots, 3\}$  such that  $|\hat{\zeta}_{p_k}(\mathbf{s})| = \max\{|\hat{\zeta}_{p_{k'}}(\mathbf{s})| : 1 \leq k' \leq 3\}$ . And for this  $k$ , we have that  $|\hat{\zeta}_{p_k}(\mathbf{s})| / \|\hat{\zeta}_p(\mathbf{s})\| \geq 1/\sqrt{3} > 0$ .

Now since  $hn^{-1/2}b_n \rightarrow \infty$ , we have that  $hn^{-1/2} \phi_p(\mathbf{s}) \hat{\zeta}_p(\mathbf{s}) / \|\hat{\zeta}_p(\mathbf{s})\| \geq hb_n / \sqrt{3n} \rightarrow \infty$  and therefore the term to the left of the equals sign dominates the sum to the right of the equals sign in (12). So for large enough  $n$ ,  $\hat{\zeta}_p(\mathbf{s}) \neq 0$  cannot maximize  $Q$ .

So  $P \left\{ \hat{\zeta}_{(b)}(\mathbf{s}) = 0 \right\} \rightarrow 1$ .

## References