

# LAGR and its oracle properties

Wesley Brooks

June 22, 2014

## 1 Introduction

Here goes the literature review. Essential elements are:

- nonparametric regression
- varying coefficients regression
- adaptive lasso
- generalized linear models

## 2 Varying coefficient models

### 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  $D \subset \mathbb{R}^2$  according

to a density  $f(\mathbf{s})$ . 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 model where the coefficients  $\boldsymbol{\beta}(\mathbf{s}_i)$  may be spatially-varying 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,

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

## 2.2 Augment the covariates and the coefficients with location interactions

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 Fan and Gijbels (1996). Here, locally linear coefficients are estimated by augmenting the local design matrix with covariate-by-location interactions in two dimensions as proposed by Wang, Mei, and Yan (2008). The augmented local design matrix at location  $\mathbf{s}_i$  is

$$\mathbf{Z}(\mathbf{s}_i) = (\mathbf{X} \quad L_i \mathbf{X} \quad M_i \mathbf{X}) \quad (3)$$

where  $\mathbf{X}$  is the unaugmented matrix of covariates,  $L_i = \text{diag}\{s_{i_1'} - s_{i_1}\}$  and  $M_i = \text{diag}\{s_{i_2'} - s_{i_2}\}$  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  $\nabla_v$ :

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

### 2.3 Local likelihood

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 since 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 likelihood, which is maximized at location  $\mathbf{s}$  to estimate the local coefficients  $\boldsymbol{\zeta}(\mathbf{s})$ :

$$\mathcal{L}\{\boldsymbol{\zeta}(\mathbf{s})\} = \prod_{i=1}^n \left\{ (2\pi\sigma^2)^{-1/2} \exp \left[ -(1/2)\sigma^{-2} \{y(\mathbf{s}_i) - \mathbf{z}'(\mathbf{s}_i) \boldsymbol{\zeta}(\mathbf{s})\}^2 \right] \right\}^{K_h(\|\mathbf{s} - \mathbf{s}_i\|)}, \quad (5)$$

The weights are computed from a kernel function  $K_h(\cdot)$  such as the Epanechnikov kernel:

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

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

Thus, the local log-likelihood function is, up to an additive constant:

$$\ell\{\boldsymbol{\zeta}(\mathbf{s})\} = -(1/2) \sum_{i=1}^n K_h(\|\mathbf{s} - \mathbf{s}_i\|) \left[ \log \sigma^2 + \sigma^{-2} \{y(\mathbf{s}_i) - \mathbf{z}'(\mathbf{s}_i)\boldsymbol{\zeta}(\mathbf{s})\}^2 \right]. \quad (7)$$

## 2.4 Estimating the coefficients

Letting  $\mathbf{W}(\mathbf{s})$  be a diagonal weight matrix where  $W_{ii}(\mathbf{s}) = K_h(\|\mathbf{s} - \mathbf{s}_i\|)$ , the local likelihood is maximized by weighted least squares:

$$\begin{aligned} \mathcal{S}\{\boldsymbol{\zeta}(\mathbf{s})\} &= (1/2) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\}^T \mathbf{W}(\mathbf{s}) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\}^T \\ \therefore \tilde{\boldsymbol{\zeta}}(\mathbf{s}) &= \{\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}(\mathbf{s})\}^{-1} \mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Y} \end{aligned} \quad (8)$$

Now Theorem 3 of Sun, Yan, Zhang, and Lu (2014) says that, for any given  $\mathbf{s}$

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

### 3 Local variable selection with LAGR

Estimating the local coefficients by (8) relies on *a priori* variable selection. The goal of local adaptive grouped regularization (LAGR) is to simultaneously select the locally relevant predictors and estimate the local coefficients. The proposed LAGR penalty is an adaptive  $\ell_1$  penalty akin to the adaptive group lasso Wang and Leng (2008); Zou (2006).

#### 3.1 Variable groupings

Each raw covariate in a LAGR model 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.

#### 3.2 Write down the LAGR-penalized local likelihood

The objective function for the LAGR at location  $\mathbf{s}$  is the penalized local sum of squares:

$$\begin{aligned} Q\{\zeta(\mathbf{s})\} &= \mathcal{S}\{\zeta(\mathbf{s})\} + \mathcal{J}\{\zeta(\mathbf{s})\} \\ &= (1/2) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\zeta(\mathbf{s})\}^T \mathbf{W}(\mathbf{s}) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\zeta(\mathbf{s})\} + \sum_{j=1}^p \phi_j(\mathbf{s}) \|\zeta_j(\mathbf{s})\| \end{aligned} \tag{9}$$

which is the sum of the weighted sum of squares  $\mathcal{S}\{\zeta(\mathbf{s})\}$  and the LAGR penalty  $\mathcal{J}\{\zeta(\mathbf{s})\}$ .

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 the local tuning parameter applied to all coefficients at location  $\mathbf{s}$  and  $\tilde{\zeta}_j(\mathbf{s})$  is the vector of unpenalized local coefficients from (8).

### 3.3 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*

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

**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  $P \left\{ \|\hat{\zeta}_j(\mathbf{s})\| = 0 \right\} \rightarrow 0$  if  $j \leq p_0$  and  $P \left\{ \|\hat{\zeta}_j(\mathbf{s})\| = 0 \right\} \rightarrow 1$  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 Sun et al. (2014), 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  $hn^{-1/2}\lambda_n(\mathbf{s})\|\zeta_j(\mathbf{s})\|^{-\gamma} \rightarrow \infty$  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  $\|\tilde{\zeta}_p(\mathbf{s})\| = O(h^{-1}n^{-1/2})$ . Then  $h^{-1}n^{-1/2}\lambda_n(\mathbf{s}) = O(n^{-1/3+\alpha})$  and  $hn^{-1/2}\lambda_n(\mathbf{s})\|\tilde{\zeta}_p(\mathbf{s})\|^{-\gamma} = O(n^{-2/3+\alpha+\gamma/3})$ .

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

## 4 Extension to GLLMs

### 4.1 Model

Generalized linear models (GLM) extend the linear model to distributions other than gaussian. The generalized local linear model (GLLM) is an extension of the GLM to varying coefficient models via local regression.

As was the case for local linear regression models, the GLLM coefficients are smooth functions of location, called  $\beta(\mathbf{s})$ . If the response variable  $y$  is from an exponential-family distribution then its density is

$$f\{y(\mathbf{s})|\mathbf{x}(\mathbf{s})\} = c\{y(\mathbf{s})\} \times \exp\left[\frac{\theta(\mathbf{s})y(\mathbf{s}) - b\{\theta(\mathbf{s})\}}{a\{\phi(\mathbf{s})\}}\right]$$

where  $\phi$  and  $\theta$  are parameters and

$$\begin{aligned} E\{y(\mathbf{s})|\mathbf{x}(\mathbf{s})\} &= \mu(\mathbf{s}) = b'\{\theta(\mathbf{s})\} \\ \theta(\mathbf{s}) &= (g \circ b')^{-1}\{\eta(\mathbf{s})\} \\ \eta(\mathbf{s}) &= \mathbf{x}^T(\mathbf{s})\beta(\mathbf{s}) = g\{\mu(\mathbf{s})\} \\ \text{Var}\{y(\mathbf{s})|\mathbf{x}(\mathbf{s})\} &= b''\{\theta(\mathbf{s})\} a\{\phi(\mathbf{s})\} \end{aligned}$$

The function  $g(\cdot)$  is called the link function. If its inverse  $g^{-1}(\cdot) = b'(\cdot)$  then the composition  $(g \circ b')(\cdot)$  is the identity function. This particular choice of  $g$

is called the canonical link. We follow the practice of ? in assuming the use of the canonical link because it is simple and because using an alternative link function would hardly affect the local fit.

Under the canonical link function, the expressions for the mean and variance of the response variable can be simplified to

$$\begin{aligned} E\{y(\mathbf{s})|\mathbf{x}(\mathbf{s})\} &= g^{-1}\{\eta(\mathbf{s})\} \\ \text{Var}\{y(\mathbf{s})|\mathbf{x}(\mathbf{s})\} &= a\{\phi(\mathbf{s})\}/g'\{\mu(\mathbf{s})\} = V\{\mu(\mathbf{s})\} \times a\{\phi(\mathbf{s})\} \end{aligned}$$

## 4.2 Local quasi-likelihood

Assuming the canonical link, all that is required is to specify the mean-variance relationship via the variance function,  $V\{\mu(\mathbf{s})\}$ . Then the GLLM coefficients can be estimated by maximizing the local quasi-likelihood

$$\ell\{\boldsymbol{\zeta}(\mathbf{s})\} = \sum_{i=1}^n K_h(\|\mathbf{s} - \mathbf{s}_i\|) Q\left[g^{-1}\{z'(\mathbf{s}_i)\boldsymbol{\zeta}(\mathbf{s})\}, Y(\mathbf{s}_i)\right]. \quad (10)$$

The local quasi-likelihood generalizes the local log-likelihood that was used to estimate coefficients in the local linear model case. The quasi-likelihood is convex, and is defined in terms of its derivative, the quasi-score function

$$\frac{\partial}{\partial \mu} Q(\mu, y) = \frac{y - \mu}{V(\mu)}.$$



### 4.3 Estimation

Under these conditions, the local quasi-likelihood is maximized where

$$\ell' \left\{ \hat{\boldsymbol{\zeta}}(\mathbf{s}) \right\} = \sum_{i=1}^n K_h(\|\mathbf{s} - \mathbf{s}_i\|) \frac{y(\mathbf{s}_i) - \hat{\mu}(\mathbf{s}_i)}{V\{\hat{\mu}(\mathbf{s}_i)\} g'\{\hat{\mu}(\mathbf{s}_i)\}} \mathbf{z}(\mathbf{s}_i) = \mathbf{0}_{3p}. \quad (11)$$

Except for the  $K_h(\|\mathbf{s} - \mathbf{s}_i\|)$  term, this is the same as the normal equations for estimating coefficients in a GLM. The method of iteratively reweighted least squares (IRLS) is used to solve for  $\hat{\boldsymbol{\zeta}}(\mathbf{s})$ .

### 4.4 Distribution of the local coefficients

The asymptotic distribution of the local coefficients in a varying-coefficients GLM with a one-dimensional effect-modifying parameter are given in ?. For coefficients that vary in two dimensions (e.g. spatial location), the asymptotic distribution under the canonical link is

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

### 4.5 LAGR penalty

As in the case of linear models, the LAGR for GLMs is a grouped  $\ell_1$  regularization method.

## 4.6 Oracle properties of LAGR in the GLM setting

## 5 Simulations

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.

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  $x_1$ , was nonzero in at least part of the domain, with a minimum of zero and maximum of one. Three parameters were varied to produce the twelve settings: there were

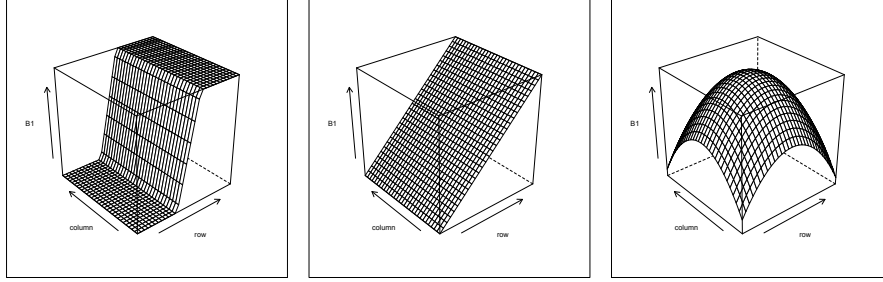


Figure 1: These three functions were used to simulate the coefficient of the covariate  $X_1$  in the VCR model  $y(\mathbf{s}_i) = x_1(\mathbf{s}_i)\beta_1(\mathbf{s}_i) + \varepsilon(\mathbf{s}_i)$ .

three functional forms for the coefficient surface  $\beta_1(\mathbf{s})$ , 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 error term. The twelve simulation settings are described in Table \ref{table:simulation\_settings}.

The three coefficient surfaces used to produce the response variable in the simulations are pictured in Figure \ref{fig:sim-actual}. 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 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.

The performance of LAGR was compared to that of a VCR model without any variable selection and a VCR model with oracular selection. Oracular selection means that exactly the correct set of covariates was used to fit each local model.

## 6 Data example

Use LAGR to estimate coefficients for a real data example. Discuss the results, perhaps in comparison to some existing analysis.

## A Proofs of theorems

*Proof of theorem 1.*

□

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

$$\begin{aligned}
V_4^{(n)}(\mathbf{u}) &= Q \left\{ \boldsymbol{\zeta}(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u} \right\} - Q \left\{ \boldsymbol{\zeta}(\mathbf{s}) \right\} \\
&= (1/2) \left[ \mathbf{Y} - \mathbf{Z}(\mathbf{s}) \left\{ \boldsymbol{\zeta}(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u} \right\} \right]^T \mathbf{W}(\mathbf{s}) \left[ \mathbf{Y} - \mathbf{Z}(\mathbf{s}) \left\{ \boldsymbol{\zeta}(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u} \right\} \right] \\
&\quad + \sum_{j=1}^p \phi_j(\mathbf{s}) \|\boldsymbol{\zeta}_j(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u}_j\| \\
&\quad - (1/2) \left\{ \mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s}) \right\}^T \mathbf{W}(\mathbf{s}) \left\{ \mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s}) \right\} - \sum_{j=1}^p \phi_j(\mathbf{s}) \|\boldsymbol{\zeta}_j(\mathbf{s})\| \\
&= (1/2) \mathbf{u}^T \left\{ h^{-2}n^{-1} \mathbf{Z}^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \mathbf{Z}(\mathbf{s}) \right\} \mathbf{u} - \mathbf{u}^T \left[ h^{-1}n^{-1/2} \mathbf{Z}^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \left\{ \mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s}) \right\} \right] \\
&\quad + \sum_{j=1}^p n^{-1/2} \phi_j(\mathbf{s}) n^{1/2} \left\{ \|\boldsymbol{\zeta}_j(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u}_j\| - \|\boldsymbol{\zeta}_j(\mathbf{s})\| \right\} \tag{12}
\end{aligned}$$

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  $\phi_j(\mathbf{s}) (\|\boldsymbol{\zeta}_j(\mathbf{s}) + h^{-1}n^{-1/2}\mathbf{u}_j\| - \|\boldsymbol{\zeta}_j(\mathbf{s})\|) = \phi_j(\mathbf{s})h^{-1}n^{-1/2}\|\mathbf{u}_j\|$ .

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

$$V_4^{*(n)}(\mathbf{u}) = \begin{cases} (1/2)\mathbf{u}^T \{h^{-2}n^{-1}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}(\mathbf{s})\} \mathbf{u} - \mathbf{u}^T [h^{-1}n^{-1/2}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\}] & \text{if } \|\mathbf{u}_j\| = 0 \forall j > p_0 \\ \infty & \text{otherwise} \end{cases}$$

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

$$\begin{aligned} 0 &= \{h^{-2}n^{-1}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}(\mathbf{s})\} \hat{\mathbf{u}}^{(n)} - [h^{-1}n^{-1/2}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\}] \\ \therefore \hat{\mathbf{u}}^{(n)} &= \{n^{-1}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}(\mathbf{s})\}^{-1} [hn^{-1/2}\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\}] \end{aligned} \quad (13)$$

By the epiconvergence results of Geyer (1994) and Knight and Fu (2000), the minimizer of the limiting function is the limit of the minimizers  $\hat{\mathbf{u}}^{(n)}$ . And since, by Lemma 2 of Sun et al. (2014),

$$\hat{\mathbf{u}}^{(n)} \xrightarrow{d} N\left(\frac{\kappa_2 h^2}{2\kappa_0} \{\nabla_{uu}^2 \boldsymbol{\zeta}_j(\mathbf{s}) + \nabla_{vv}^2 \boldsymbol{\zeta}_j(\mathbf{s})\}, f(\mathbf{s})\kappa_0^{-2}\nu_0\sigma^2\Psi^{-1}\right) \quad (14)$$

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  $P\left\{\hat{\zeta}_j(\mathbf{s}) = 0\right\} \rightarrow 1$  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

$$\begin{aligned}
0 &= \mathbf{Z}_p^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \left\{ \mathbf{Y} - \mathbf{Z}_{-p}(\mathbf{s}) \hat{\zeta}_{-p}(\mathbf{s}) - \mathbf{Z}_p(\mathbf{s}) \hat{\zeta}_p(\mathbf{s}) \right\} - \phi_p(\mathbf{s}) \frac{\hat{\zeta}_p(\mathbf{s})}{\|\hat{\zeta}_p(\mathbf{s})\|} \\
&= \mathbf{Z}_p^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \left[ \mathbf{Y} - \mathbf{Z}(\mathbf{s}) \zeta(\mathbf{s}) - \frac{h^2 \kappa_2}{2\kappa_0} \left\{ \nabla_{uu}^2 \zeta(\mathbf{s}) + \nabla_{vv}^2 \zeta(\mathbf{s}) \right\} \right] \\
&\quad + \mathbf{Z}_p^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \mathbf{Z}_{-p}(\mathbf{s}) \left[ \zeta_{-p}(\mathbf{s}) + \frac{h^2 \kappa_2}{2\kappa_0} \left\{ \nabla_{uu}^2 \zeta_{-p}(\mathbf{s}) + \nabla_{vv}^2 \zeta_{-p}(\mathbf{s}) \right\} - \hat{\zeta}_{-p}(\mathbf{s}) \right] \\
&\quad + \mathbf{Z}_p^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \mathbf{Z}_p(\mathbf{s}) \left[ \zeta_p(\mathbf{s}) + \frac{h^2 \kappa_2}{2\kappa_0} \left\{ \nabla_{uu}^2 \zeta_p(\mathbf{s}) + \nabla_{vv}^2 \zeta_p(\mathbf{s}) \right\} - \hat{\zeta}_p(\mathbf{s}) \right] \\
&\quad - \phi_p(\mathbf{s}) \frac{\hat{\zeta}_p(\mathbf{s})}{\|\hat{\zeta}_p(\mathbf{s})\|}
\end{aligned} \tag{15}$$

So

$$\begin{aligned}
\frac{h}{\sqrt{n}} \phi_p(\mathbf{s}) \frac{\hat{\zeta}_p(\mathbf{s})}{\|\hat{\zeta}_p(\mathbf{s})\|} &= \mathbf{Z}_p^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \frac{h}{\sqrt{n}} \left[ \mathbf{Y} - \mathbf{Z}(\mathbf{s}) \zeta(\mathbf{s}) - \frac{h^2 \kappa_2}{2\kappa_0} \left\{ \nabla_{uu}^2 \zeta(\mathbf{s}) + \nabla_{vv}^2 \zeta(\mathbf{s}) \right\} \right] \\
&\quad + \left\{ n^{-1} \mathbf{Z}_p^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \mathbf{Z}_{-p}(\mathbf{s}) \right\} h \sqrt{n} \left[ \zeta_{-p}(\mathbf{s}) + \frac{h^2 \kappa_2}{2\kappa_0} \left\{ \nabla_{uu}^2 \zeta_{-p}(\mathbf{s}) + \nabla_{vv}^2 \zeta_{-p}(\mathbf{s}) \right\} - \hat{\zeta}_{-p}(\mathbf{s}) \right] \\
&\quad + \left\{ n^{-1} \mathbf{Z}_p^T(\mathbf{s}) \mathbf{W}(\mathbf{s}) \mathbf{Z}_p(\mathbf{s}) \right\} h \sqrt{n} \left[ \zeta_p(\mathbf{s}) + \frac{h^2 \kappa_2}{2\kappa_0} \left\{ \nabla_{uu}^2 \zeta_p(\mathbf{s}) + \nabla_{vv}^2 \zeta_p(\mathbf{s}) \right\} - \hat{\zeta}_p(\mathbf{s}) \right]
\end{aligned} \tag{16}$$

From Lemma 2 of Sun et al. (2014),  $\{n^{-1}\mathbf{Z}_p^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}_{-p}(\mathbf{s})\} = O_p(1)$  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 Sun et al. (2014), 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 (16) are  $O_p(1)$ .

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

The three terms of the sum to the right of the equals sign in (16) 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 (16). 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

Cai, Z., J. Fan, and R. Li (2000). Efficient estimation and inferences for varying-coefficient models. *Journal of the American Statistical Association* 95(888-902).

- Fan, J. and I. Gijbels (1996). *Local Polynomial Modeling and its Applications*. Chapman and Hall, London.
- Fan, J., N. E. Heckman, and M. Wand (1995). Local polynomial kernel regression for generalized linear models and quasi-likelihood functions. *Journal of the American Statistical Association* 90, 141–150.
- Geyer, C. J. (1994). On the asymptotics of constrained m-estimation. *Annals of Statistics* 22, 1993–2010.
- Knight, K. and W. Fu (2000). Asymptotics for lasso-type estimators. *Annals of Statistics* 28, 1356–1378.
- Sun, Y., H. Yan, W. Zhang, and Z. Lu (2014). A semiparametric spatial dynamic model. *Annals of Statistics*.
- Wang, H. and C. Leng (2008). A note on adaptive group lasso. *Computational Statistics and Data Analysis* 52, 5277–5286.
- 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.
- Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* 101, 1418–1429.