

# Oracle properties of local adaptive grouped regularization

Wesley Brooks

---

---

## 1. Spatially varying coefficients regression

### 1.1. 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 observation location  $\mathbf{s}_i$ , assume that the outcome is related to the covariates by a generalized linear model (GLM) where the coefficients  $\boldsymbol{\beta}(\mathbf{s}_i)$  may be spatially varying. That is, the distribution of  $Y$  conditional on  $X$  is:

$$f_{Y|X}(y|x) = \exp \{y\}$$

$$\mu(\mathbf{s}_i) = g^{-1} \{ \eta(\mathbf{s}_i) \} \tag{1}$$

$$\eta(\mathbf{s}_i) = \mathbf{x}(\mathbf{s}_i)' \boldsymbol{\beta}(\mathbf{s}_i) \tag{2}$$

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,

$$\boldsymbol{\varepsilon} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2). \tag{3}$$

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, locally linear coefficients are estimated by augmenting the local design matrix with covariate-by-location interactions in two dimensions as proposed by ?. The augmented local design matrix at location  $\mathbf{s}_i$  is

$$\mathbf{Z}(\mathbf{s}_i) = (\mathbf{X} \ L_i \mathbf{X} \ M_i \mathbf{X}) \quad (4)$$

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  $\nabla_v$ :

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

## 1.2. Estimation

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

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

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  $\boldsymbol{\zeta}(\mathbf{s})$ :

$$\mathcal{L}\{\boldsymbol{\zeta}(\mathbf{s})\} = \prod_{i=1}^n \left\{ (2\pi\sigma^2)^{-1/2} \exp \left[ -(1/2)\sigma^{-2} \{y(\mathbf{s}_i) - \mathbf{z}'(\mathbf{s}_i)\boldsymbol{\zeta}(\mathbf{s})\}^2 \right] \right\}^{K_h(\|\mathbf{s}-\mathbf{s}_i\|)}, \quad (6)$$

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

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

$$K(x) = \begin{cases} (3/4)(1-x^2) & \text{if } x < 1, \\ 0 & \text{if } x \geq 1. \end{cases} \quad (7)$$

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

$$\ell\{\boldsymbol{\zeta}(\mathbf{s})\} = -(1/2) \sum_{i=1}^n K_h(\|\mathbf{s} - \mathbf{s}_i\|) \left[ \log \sigma^2 + \sigma^{-2} \{y(\mathbf{s}_i) - \mathbf{z}'(\mathbf{s}_i)\boldsymbol{\zeta}(\mathbf{s})\}^2 \right]. \quad (8)$$

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:

$$\begin{aligned} \mathcal{S}\{\boldsymbol{\zeta}(\mathbf{s})\} &= (1/2) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\}^T \mathbf{W}(\mathbf{s}) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\}^T \\ \therefore \tilde{\boldsymbol{\zeta}}(\mathbf{s}) &= \{\mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Z}(\mathbf{s})\}^{-1} \mathbf{Z}^T(\mathbf{s})\mathbf{W}(\mathbf{s})\mathbf{Y} \end{aligned} \quad (9)$$

## 2. Local variable selection and parameter estimation

### 2.1. Local variable selection

Local adaptive grouped regularization (LAGR) is explored as a method of local variable selection and coefficient estimation in SVCR models. The proposed LAGR penalty is an adaptive  $\ell_1$  penalty

akin to the adaptive group lasso (??).

Grouped variables are selected together for inclusion in the model. Each group in a LAGR model consists of one covariate and its gradients on the two dimensions of spatial location. That is,

$$\boldsymbol{\zeta}_j(\mathbf{s}) = (\beta_j(\mathbf{s}) \quad \nabla_u \beta_j(\mathbf{s}) \quad \nabla_v \beta_j(\mathbf{s}))^T \text{ for } j = 1, \dots, p.$$

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

$$\begin{aligned} Q\{\boldsymbol{\zeta}(\mathbf{s})\} &= \mathcal{S}\{\boldsymbol{\zeta}(\mathbf{s})\} + \mathcal{J}\{\boldsymbol{\zeta}(\mathbf{s})\} \\ &= (1/2) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\}^T \mathbf{W}(\mathbf{s}) \{\mathbf{Y} - \mathbf{Z}(\mathbf{s})\boldsymbol{\zeta}(\mathbf{s})\} + \sum_{j=1}^p \phi_j(\mathbf{s}) \|\boldsymbol{\zeta}_j(\mathbf{s})\| \end{aligned} \quad (10)$$

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

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

## 2.2. Computation

### 2.2.1. Tuning parameter selection

Implementing LAGR requires the selection of local tuning parameters. The criteria commonly used for selecting tuning parameters in lasso-type models are appropriate here, including GCV (?), Cp (?), AIC (?), and BIC (?). All of the examples and simulations presented herein used the corrected AIC (AICc) (?) for tuning parameter selection:

$$\text{AIC}_c(\mathbf{s}) = \hat{\sigma}^{-2}(\mathbf{s}) \|\mathbf{Y} - \mathbf{Z}(\mathbf{s})\hat{\boldsymbol{\zeta}}(\mathbf{s})\|^2 + 2df + 2 \frac{df(df+1)}{\sum_{i=1}^n K_h(\|\mathbf{s} - \mathbf{s}_i\|) - df - 1} \quad (11)$$

where  $df$  is the degrees of freedom as defined in ?:

$$df = \sum_{j=1}^p I \left\{ \|\hat{\zeta}_j(\mathbf{s})\| \right\} + 2 \sum_{j=1}^p \frac{\|\hat{\zeta}_j(\mathbf{s})\|}{\|\tilde{\zeta}_j(\mathbf{s})\|} \quad (12)$$

### 3. References