

Spatial Regression Using Kernel Averaged Predictors

Matthew J. HEATON and Alan E. GELFAND

Traditional spatial linear regression models assume that the mean of the response is a linear combination of predictors measured at the same location as the response. In spatial applications, however, it seems plausible that neighboring predictors can also inform about the response. This article proposes using unobserved kernel averaged predictors in such regressions. The kernels are parametric introducing additional parameters that are estimated with the data. Properties and challenges of using kernel averaged predictors within a regression model are detailed in the simple case of a univariate response and a single predictor. Additionally, extensions to multiple predictors and generalized linear models are discussed. The methods are demonstrated using a data set of dew duration and shrub density. Supplemental materials are available online.

Key Words: Block averaging; Circular neighborhoods; Distributed covariates; Multivariate Gaussian process; Spatial linear model.

1. INTRODUCTION

Statistical methods and models for spatial data are plentiful, and panoramic overviews of spatial methodology can be found in the books by Cressie (1993), Chiles and Delfiner (1999), Wackernagel (2003), and Banerjee, Carlin, and Gelfand (2004). In spatial problems, a response variable is measured at various locations over a region \mathcal{D} . The response variable is correlated in space such that observations closer in space tend to be more correlated than observations farther away in space. For continuous responses, a normality assumption (perhaps on a transformation of the response variable) is typically used, and spatial correlation is induced via Gaussian processes. For binary or count data, spatial correlation is typically induced via random effects on the mean function as in Diggle, Tawn, and Moyeed (1998). This article primarily focuses on continuous response data, but extensions to binary or count data are discussed in Section 5.

Matthew J. Heaton is Ph.D. Student (E-mail: matt@stat.duke.edu) and Alan E. Gelfand (✉) is Professor (E-mail: alan@stat.duke.edu), Department of Statistical Science, Duke University, Box 90251, Durham, NC 27708-0251, USA

© 2011 International Biometric Society
Journal of Agricultural, Biological, and Environmental Statistics, Volume 16, Number 2, Pages 233–252
DOI: [10.1007/s13253-010-0050-6](https://doi.org/10.1007/s13253-010-0050-6)

In many spatial applications, a vector of predictors (covariates) is also available, and the goal is to build a regression model for the response while accounting for spatial correlation in the response, the predictors, or both. For such spatial regression models, predictors are most commonly incorporated linearly via the mean function such that $\mathbb{E}(Y(s) \mid \beta_0, \boldsymbol{\beta}, \mathbf{X}(s)) = \beta_0 + \boldsymbol{\beta}^T \mathbf{X}(s)$, where $Y(s)$ is a univariate response at location $s \in \mathcal{D}$, $\mathbf{X}(s) = (X_1(s), \dots, X_p(s))^T$ is a vector of p predictor variables measured at location s , and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$ is the vector of coefficients. An extension of the spatial regression model introduces local linearity through spatially varying coefficients as in Gelfand et al. (2003), who assume that $\mathbb{E}(Y(s) \mid \boldsymbol{\beta}(s), \mathbf{X}(s)) = \beta_0(s) + \boldsymbol{\beta}^T(s) \mathbf{X}(s)$ so that the coefficients are different for all s .

In formulating a spatial regression model it is often necessary to assume that $\mathbf{X}(s)$ is random. This enables handling of the full dataset including observations where an X or Y is missing and hence handling *misalignment* in the dataset (see Banerjee, Carlin, and Gelfand 2004, Chapter 6). It also enables handling of the inverse problem, learning about $\mathbf{X}(s_0)$ for an observed $Y(s_0)$. This implies a joint model for $\mathbf{Z}(s) = (Y(s), \mathbf{X}^T(s))^T$ as developed, for instance, in Royle and Berliner (1999), Berliner (2000), or Gelfand et al. (2004). Analogous to the case where the $\mathbf{X}(s)$ are viewed as fixed, the local spatial regression is usually defined in terms of the conditional distribution of $Y(s)$ given $\mathbf{X}(s)$ (see Banerjee, Carlin, and Gelfand 2004, Chapter 7) in the form of the conditional mean $E(Y(s) \mid \mathbf{X}(s)) = \beta_0 + \boldsymbol{\beta}^T \mathbf{X}(s)$, as above. However, in spatial applications, *neighboring* $\mathbf{X}(s')$ can be expected to inform about $Y(s)$ particularly when the distance between s and s' is small. To cite a few examples, precipitation can affect the water table (the depth at which soil and pore spaces become completely saturated with water) at not only where the precipitation fell but also surrounding locations due to run off and changes in slope from uneven ground surfaces. Similarly, the concentration of ozone (O_3) is affected by pollutants, ultraviolet rays, and temperature within a neighborhood of location s . Thus, a mean specification which only includes $\mathbf{X}(s)$ and not neighboring $\mathbf{X}(s')$ may not adequately capture the process.

A mean specification which incorporates neighboring predictors might take the form $\mathbb{E}(Y(s) \mid \boldsymbol{\theta}, \mathbf{X}) = \beta_0(s) + \int_{\mathcal{D}} \boldsymbol{\beta}(s, \mathbf{u} \mid \boldsymbol{\theta})^T \mathbf{X}(\mathbf{u}) d\mathbf{u}$, where $\boldsymbol{\beta}(s, s' \mid \boldsymbol{\theta})$ is a coefficient model capturing the effect of $\mathbf{X}(s')$ on $Y(s)$ with parameters $\boldsymbol{\theta}$. In the interest of clarification, with, say n observations $\mathbf{Z} = (\mathbf{Z}(s_1), \dots, \mathbf{Z}(s_n))^T$, we might consider the conditional distribution, $[Y(s_i) \mid \mathbf{X}(s_1), \dots, \mathbf{X}(s_n)]$, where $[\cdot]$ denotes a generic probability distribution. In the Gaussian setting, $[Y(s_i) \mid \mathbf{X}(s_1), \dots, \mathbf{X}(s_n)]$ would contain n coefficients $\boldsymbol{\beta}(s_i, s_1), \dots, \boldsymbol{\beta}(s_i, s_n)$ which describe the loading of $\mathbf{X}(s_j)$ on $Y(s_i)$ for $j = 1, \dots, n$. Such a choice is analogous to spatial prediction or kriging and is too limited for our goal because the explained relationship between the response and predictor is dependent on the number of sampling locations (n) and the arrangement of those locations in \mathcal{D} . It also fails to capture the idea of local $\mathbf{X}(s')$ informing about $Y(s)$.

Below, we argue that kernels are an appropriate choice for $\boldsymbol{\beta}(s, s' \mid \boldsymbol{\theta})$ and present associated models. We offer full distributional detail and properties, show how to fit these models, and illustrate benefits using real and simulated data. We close this section with three remarks. First, we are not interested in what would be called co-kriging, i.e., prediction of $Y(s_0)$ given the entire set of observed X 's and Y 's. When used for spatial prediction, the models considered here offer only minimal additional benefit from what would be

obtained using customary multivariate Gaussian processes. This is not surprising. We are introducing unobserved covariates into our regression specification, but we are still fitting a multivariate Gaussian process model to the data. We offer further discussion of this point in Section 2.2 below. Second, what we are doing differs from what is usually referred to as employing functional covariates (see Baíllo and Grané 2009, and the references therein). Functional covariates envision a function at location s , say, $X(s, t)$ for $t \in (0, T]$, and seek to reduce this to a single covariate at s to explain $Y(s)$. Usually, this is achieved through some integration of the function wherein the integration is over t rather than over s as we seek. Third, we do not seek to build process models for $\beta(s, s' | \theta)$. We provide process modeling for $(Y(s), X(s))$ but specify β as a parametric function. Specifying the latter using processes will provide poorly-identified, over-fitted models.

Section 2 develops the methodology for the simplest case of a univariate response and a single predictor, and Section 3 presents the associated model fitting. Section 4 uses this methodology on simulated and field data collected over the Negev Desert in Israel. Section 5 discusses extensions of the methodology to multiple predictors and general (e.g., count or binary) response variables. Conclusions and future work are discussed in Section 6.

2. SPATIAL REGRESSION WITH A SINGLE AVERAGED COVARIATE

2.1. THE MODEL

Let $Y(s)$ and $X(s)$ denote a univariate response variable and a single covariate at location $s \in \mathcal{D} \subseteq \mathbb{R}^d$ for $d \in \mathbb{Z}$, respectively. Furthermore, assume that $X(s)$ follows a Gaussian process (GP) of the form

$$X(s) = \mu_X(s) + \sigma_X w_X(s), \quad (2.1)$$

where $\mu_X(s)$ is the mean surface at location s , and $w_X(s)$ is a mean 0, variance 1 GP with correlation function $\text{corr}(w_X(s), w_X(s')) = \rho_X(s, s' | \phi_X)$, where ϕ_X denotes the (possibly vector-valued) parameter associated with ρ_X . For example, if ρ_X is the Matérn (Matérn 1986) correlation function, then $\phi_X = (\psi_X, \nu_X)^T$, where ψ_X is the decay parameter, and ν_X is the smoothness parameter. Notice that (2.1) defines a purely spatial covariate process. In some applications, however, covariates are measured with error. In these cases, $X(s)$ can be thought of as the “true” underlying covariate process, and the observed covariate is $H(s) = X(s) + \epsilon_X(s)$, where $\epsilon_X(s)$ is an *iid* white noise process. Distributional results for $H(s)$ will differ from the those of $X(s)$ by an additive nugget variance term only, so, for simplicity, we assume that $X(s)$ is observed directly for the remainder of this article. Moreover, we still want to use $X(s)$ to explain $Y(s)$, i.e., we want the regression of $Y(s)$ to be on the true $X(s)$.

Following the discussion in the Introduction, we define the unobserved local covariate at s incorporating neighbor information as $\tilde{X}(s)$ using a kernel function, i.e.,

$$\tilde{X}(s) \equiv \frac{1}{K(s | \xi)} \int_{\mathcal{D}} K(s, u | \xi) X(u) du, \quad (2.2)$$

where $K(s, s' | \xi)$ is a kernel defining a weight on the distance between s and s' with parameters ξ and $0 < K(s | \xi) = \int_{\mathcal{D}} K(s, u | \xi) du < \infty$. The parameter ξ , most commonly, consists of *scale* parameters such as the entries of a covariance matrix but can also include *location* parameters as in Higdon (1998) and Xu, Wikle, and Fox (2005). Choices for K are discussed below. By (2.2), $\tilde{X}(s)$ is the kernel averaged value of $X(s)$ over the domain \mathcal{D} . Again, we note that the integration is over \mathcal{D} , i.e., $X(s)$ is not viewed as a functional covariate. The latter would take the form, say $X(s, t)$, whence the integration would be over t .

Because a valid GP was defined for $X(s)$, $\tilde{X}(s)$ also is a valid GP with mean $\tilde{\mu}_X(s) = K(s | \xi)^{-1} \int_{\mathcal{D}} K(s, u | \xi) \mu_X(u) du$ and

$$\begin{aligned} \text{Cov}(\tilde{X}(s), \tilde{X}(s')) &= \frac{\sigma_X^2}{K(s | \xi) K(s' | \xi)} \int_{\mathcal{D}} \int_{\mathcal{D}} K(s, u | \xi) K(s', v | \xi) \rho_X(u, v) dv du \\ &\equiv \sigma_X^2 \rho_{\tilde{X}}(s, s'). \end{aligned} \quad (2.3)$$

Not only are $X(s)$ and $\tilde{X}(s)$ marginally Gaussian processes, but a valid bivariate GP is induced for the pair $(X(s), \tilde{X}(s))^T$. Specifically, $(X(s), \tilde{X}(s))^T$ follows a bivariate GP with mean $(\mu_X(s), \tilde{\mu}_X(s))^T$ and $\text{Cov}((X(s), \tilde{X}(s))^T, (X(s'), \tilde{X}(s'))^T)$ given by

$$\begin{aligned} \Sigma_{X\tilde{X}}(s, s') &= \sigma_X^2 \begin{pmatrix} \rho_X(s, s') & \frac{1}{K(s' | \xi)} \int_{\mathcal{D}} K(s', u | \xi) \rho_X(s, u) du \\ \frac{1}{K(s | \xi)} \int_{\mathcal{D}} K(s, u | \xi) \rho_X(u, s') du & \rho_{\tilde{X}}(s, s') \end{pmatrix} \\ &\equiv \sigma_X^2 \begin{pmatrix} \rho_X(s, s') & \rho_{X, \tilde{X}}(s, s') \\ \rho_{\tilde{X}, X}(s, s') & \rho_{\tilde{X}}(s, s') \end{pmatrix} \end{aligned} \quad (2.4)$$

for any other location $s' \in \mathcal{D}$.

To account for effects of $\{X(s') : s' \in \mathcal{D}\}$ on $Y(s)$, consider the linear model defined by

$$Y(s) | X(s), \tilde{X}(s) = \beta_0 + \beta_1 \tilde{X}(s) + \sigma_Y w_Y(s) + \epsilon_Y(s), \quad (2.5)$$

where $w_Y(s)$ is defined analogously to $w_X(s)$ in (2.1) but with correlation function ρ_Y , and $\epsilon_Y(s)$ is a Gaussian white noise process with variance τ_Y^2 . Intuitively, the kernel $K(s, u | \xi)$ is used here to describe how the effect of the covariate $X(s)$ propagates to the response. In contrast, the kernels used in the integro-difference equation models of Wikle (2002) and Xu, Wikle, and Fox (2005) describe how the spatio-temporal process, say $Y(s, t)$, propagates in time.

Notice that the bivariate GP for $(X(s), \tilde{X}(s))^T$ above, together with (2.5), provides a joint specification of a trivariate GP for $\mathbf{Z}(s) = (X(s), \tilde{X}(s), Y(s))^T$. Such sequential specification is discussed in Royle and Berliner (1999), Berliner (2000), and Gelfand et al. (2004). Specifically, $\mathbf{Z}(s)$ follows a valid trivariate GP with mean $\boldsymbol{\mu}(s) = (\mu_X(s), \tilde{\mu}_X(s), \beta_0 + \beta_1 \tilde{\mu}_X(s))^T$ and cross-covariance

$$\sigma_X^2 \begin{pmatrix} \rho_X(s, s') & \rho_{X, \tilde{X}}(s, s') & \beta_1 \rho_{X, \tilde{X}}(s, s') \\ \rho_{\tilde{X}, X}(s, s') & \rho_{\tilde{X}}(s, s') & \beta_1 \rho_{\tilde{X}}(s, s') \\ \beta_1 \rho_{\tilde{X}, X}(s, s') & \beta_1 \rho_{\tilde{X}}(s, s') & \frac{\sigma_Y^2}{\sigma_X^2} \rho_Y(s, s') + \beta_1^2 \rho_{\tilde{X}}(s, s') \end{pmatrix}. \quad (2.6)$$

Table 1. Examples of kernel functions $K(s, s' | \xi)$ and their parameters. $\mathbb{I}_{\{\mathcal{A}\}}$ is the indicator of a set \mathcal{A} .

Kernel	$K(s, s' \xi)$	Parameters (ξ)
Uniform	$\mathbb{I}_{\{\ s-s'\ \leq \xi\}}$	ξ
Epanechnikov	$(\xi^2 - \ s-s'\ ^2) \mathbb{I}_{\{\ s-s'\ \leq \xi\}}$	ξ
Component Wise Gaussian	$\prod_{i=1}^d \xi_i^{-1} \exp\{-(s_i - s'_i)^2 / (2\xi_i^2)\}$	$\xi_i, i = 1, \dots, d$
Oriented Gaussian	$ \Xi ^{-1/2} \exp\{-\frac{(s'-s)^T \Xi^{-1} (s'-s)}{2}\}$	$\Xi = \{\xi_{ij}\}$

The joint distribution of $\mathbf{Z}(s)$ is useful for evaluating the properties of (2.5) such as induced correlations between $Y(s)$ and $(X(s), \tilde{X}(s))$ (see Section 2.2). Note also that the induced bivariate process model for $(X(s), Y(s))$ is not a usual bivariate process specification; the kernel appears in the covariance structure.

Analogous to the discussion in Section 1, while we have a joint process model for $\mathbf{Z}(s)$, as above, we focus on the *local* spatial regression in terms of the conditional distribution of $Y(s)$ given $\tilde{X}(s)$ in the form of the conditional mean $E(Y(s)|\tilde{X}(s)) = \beta_0 + \beta_1 \tilde{X}(s)$. Thus, given $\tilde{X}(s)$, concepts such as R^2 , mean square error, variable selection, shrinkage, etc. are applicable. Of course, all are random and would be averaged over the distribution of $\tilde{X}(s)$ in order to interpret them. Evidently, the potentially complex relationship between $\{X(s) : s \in \mathcal{D}\}$ and $Y(s)$ is captured through a single parameter (β_1). However, from (2) we see that, effectively, we are introducing a coefficient weighting of the entire surface $X(s)$ to explain $Y(s)$. That is, the coefficient of $X(s')$ is $\beta_1 K(s, s' | \xi) / K(s | \xi)$. The normalization by $K(s | \xi)$ identifies β_1 . Extending β_1 to $\beta_1(s)$ using spatially varying coefficients modeling as described in Gelfand et al. (2003) could be envisioned to create an arbitrarily rich specification. However, such models are expected to be poorly identified since they introduce the product form $\beta(s)\tilde{X}(s)$ where each variable comes from a process model but neither is observed. Our proposed form through β_1 and ξ is parsimonious and well identified.

The choice of K with the data informing about ξ enables a fairly rich regression specification while attractively reducing to a simple linear regression model in $\tilde{X}(s)$. Examples of kernels are given in Table 1, and kernel selection is explored in more detail in Section 4.2. We note that computation of $K(s | \xi)$ for general regions \mathcal{D} varies with s and can be computationally expensive when done repeatedly over MCMC iterations. For much of the following, the kernel is taken to be $K(s, s' | \xi) = \mathbb{I}_{\{\|s-s'\| \leq \xi g(\phi_X)\}}$, where $\mathbb{I}_{\{\mathcal{A}\}}$ is the indicator for a set \mathcal{A} , $\|\cdot\|$ denotes Euclidean distance, $\xi \in (0, 1)$, and $g(\phi_X)$ is the effective spatial range associated with ρ_X ; for example, the exponential correlation function has $g(\phi_X) \approx 3/\phi_X$, where ϕ_X is the spatial decay parameter. Intuitively, the kernel $K(s, s' | \xi)$ is a disk centered at location s with radius $r = \xi g(\phi_X)$. The scale parameter r can be loosely interpreted as a hard threshold “decay” parameter in that the effect of $X(s')$ on $Y(s)$ is negligible if $\|s - s'\| \geq r$. Given $K(s, s' | \xi)$ as above, $K(s | \xi) = |\{s' : \|s - s'\| \leq r\} \cap \mathcal{D}|$, where $|A|$ denotes the area of A . For model fitting, if we make \mathcal{D} large enough so that each s_i in the sample and for all r of interest the disks are contained in \mathcal{D} , then $K(s_i | \xi) = \pi r^2$. This facilitates model fitting, but if we seek to

interpolate $\mathbb{E}(Y(s_0) \mid \tilde{X}(s_0))$ for s_0 near the boundary of \mathcal{D} , then we will have to use the general form $K(s \mid \xi)$. Moreover, if r is large enough so that the disks contain \mathcal{D} , then $\tilde{X}(s) = |\mathcal{D}|^{-1} \int_{\mathcal{D}} X(u) du$. In other words, for large r , $\tilde{X}(s) \approx \tilde{X}(s')$ for all s, s' yielding a highly collinear regression.

With the choice of $K(s, s' \mid \xi)$ above, r and ϕ_X are strongly associated parameters; this is evident from the forms in (2.6). For illustration, let $\rho_X(\cdot) = \exp\{-\phi_X \|s - s'\|\}$ be the exponential correlation function. Figure 1 displays contour plots of $\text{corr}(Y(s), X(s))$ and $\text{corr}(Y(s), \tilde{X}(s))$ for various values of ϕ_X and r with $\beta_1 = \sigma_X^2 = \sigma_Y^2 = 1$. In general, holding the other parameters fixed, $\text{corr}(Y(s), X(s))$ and $\text{corr}(Y(s), \tilde{X}(s))$ decrease non-linearly as r and ϕ_X increase, but the rate of decay is greater for $\text{corr}(Y(s), X(s))$. Indeed, Figure 1 suggests that plausible values for scale parameters of kernels depend on ϕ_X as well as \mathcal{D} . For the K above, parameterizing the scale parameter as $r = \xi g(\phi_X)$ where $\xi \in (0, 1)$ removes this dependency such that the prior distribution for ξ and ϕ_X can be taken as $[\xi][\phi_X]$ and restricts K to be within the effective spatial range of ρ_X . In addition, we choose $[\xi]$ and $[\phi_X]$ to be discrete to help with identifiability (see Zhang 2004) and to make computation easier (see Section 3). Similar parameterizations of scale parameters can be envisioned for other kernels (see Section 4.2).

2.2. USING $X(s)$ INSTEAD OF $\tilde{X}(s)$

We now consider what (2.5) implies about using $X(s)$ instead of $\tilde{X}(s)$ in the conditional mean for $Y(s) \mid X(s), \tilde{X}(s)$ assuming that $\mathbb{E}(Y(s) \mid X(s), \tilde{X}(s)) = \beta_0 + \beta_1 \tilde{X}(s)$. Similarly, we could investigate the consequences of assuming model (2.5) when, in fact, $\mathbb{E}(Y(s) \mid X(s), \tilde{X}(s)) = \beta_0 + \beta_1 X(s)$. We focus on the former model misspecification here, but results for the latter are given in the supplementary materials.

Using the fact that $\mathbf{Z}(s) = (X(s), \tilde{X}(s), Y(s))^T$ follows a trivariate GP with known mean and covariance given by (2.6), multivariate normal theory gives that $Y(s) \mid X(s)$ is also normally distributed with mean and variance,

$$\mathbb{E}(Y(s) \mid X(s)) = \beta_0 + \beta_1(\tilde{\mu}_X(s) - \mu_X(s)\rho_{X,\tilde{X}}(s, s)) + \beta_1\rho_{X,\tilde{X}}(s, s) \times X(s), \quad (2.7)$$

$$\text{Var}(Y(s) \mid X(s)) = \tau_Y^2 + \sigma_Y^2 + \beta_1^2\sigma_X^2(\rho_{\tilde{X}}(s, s) - \rho_{X,\tilde{X}}^2(s, s)), \quad (2.8)$$

and

$$\text{Cov}(Y(s), Y(s') \mid X(s), X(s')) = \sigma_Y^2\rho_Y(s, s') + \beta_1^2\sigma_X^2\rho_{\tilde{X}}(s, s'). \quad (2.9)$$

Notice that when the model given by (2.5) holds, then the change in $Y(s)$ as a result of a unit change in $X(s)$ is $\beta_1\rho_{X,\tilde{X}}(s, s)$ as opposed to β_1 when using $\tilde{X}(s)$. Hence, under the trivariate Gaussian process model, i.e., when the assumptions of (2.5) hold, using $X(s)$ implies that the effect of the covariate on $Y(s)$ is shrunk towards zero; $|\beta_1\rho_{X,\tilde{X}}(s, s)| \leq |\beta_1|$ because $0 \leq \rho_{X,\tilde{X}}(s, s) \leq 1$. This result is not surprising in that if other $X(s')$ in \mathcal{D} besides $X(s)$ affect $Y(s)$, then the effect due to $X(s)$ is expected to diminish. The amount of shrinkage is determined by the kernel parameters, ξ , as well as ϕ_X . Figure 2 displays $1 - \rho_{X,\tilde{X}}(s, s)$ where $\mathcal{D} = \mathbb{R}^2$, $K(s, s' \mid \xi, \phi_X) = \mathbb{I}_{\{\|s-s'\| \leq r\}}$ with $r = \xi g(\phi_X)$ and $\xi \in (0, 1)$, ρ_X is the exponential correlation function, and $\phi_X = 10$ is the spatial decay parameter.

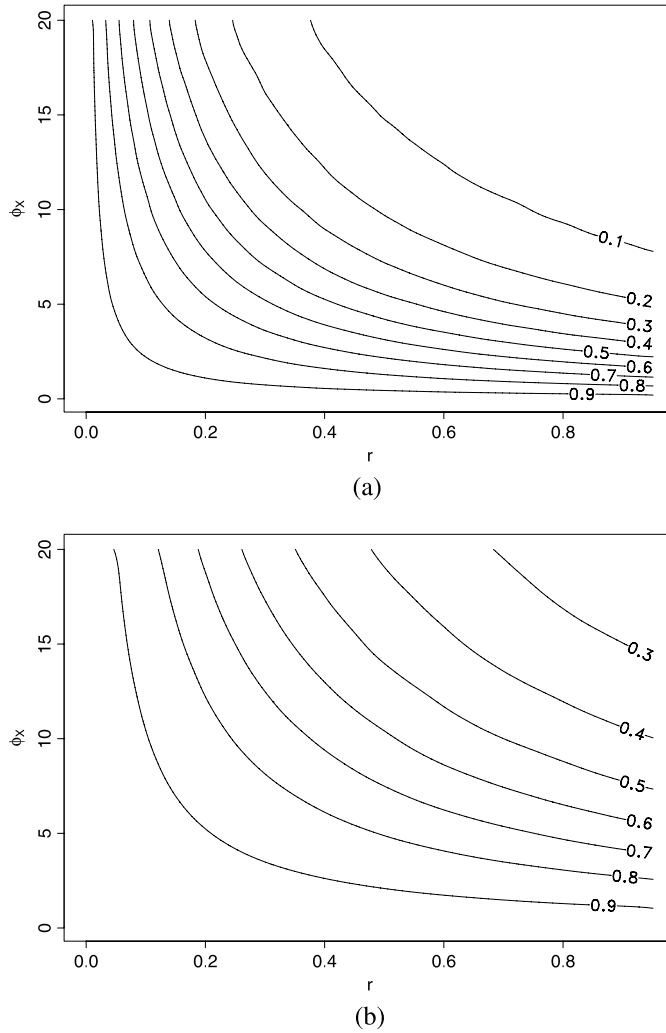


Figure 1. Contour plots of (a) $\text{corr}(Y(s), X(s))$ and (b) $\text{corr}(Y(s), \tilde{X}(s))$ as functions of r and ϕ_X . Both correlations decrease as r and ϕ_X increase, but the rate of decay is much faster for $\text{corr}(Y(s), X(s))$.

Figure 2 shows that the amount of shrinkage could be severe depending on the value of ξ . To appreciate the implications of this, suppose, for example, that $Y(s)$ is ozone and $X(s)$ is temperature (since high temperatures are known to encourage the formation of ozone). Then using $X(s)$ in the model could, potentially, underestimate the change in $Y(s)$ as a result from a unit change in temperature, leading to underestimation of the production of ozone.

A second consequence of using $X(s)$ instead of $\tilde{X}(s)$ when (2.5) holds is that the percent of variation in $Y(s)$ explained by $X(s)$ is less than the percent of variation in $Y(s)$ explained by $\tilde{X}(s)$ for many common covariance functions as detailed by the following result.

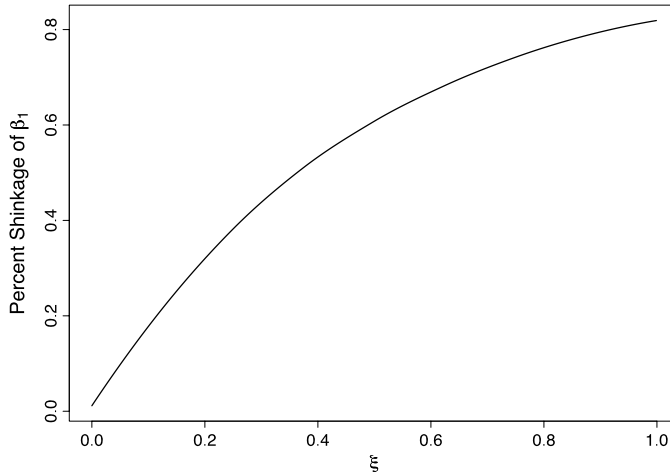


Figure 2. Percent shrinkage of β_1 against ξ assuming $\phi_X = 10$. The rapid decay of $\rho_{X,\tilde{X}}(s, s)$ as ξ increases leads to a smaller estimated effect of $X(s)$ on the response.

Result 1. Let $\rho_{Y|X}^2$ and $\rho_{Y|\tilde{X}}^2$ be the population coefficients of determination for the linear models defined by (2.7) and (2.5), respectively. If ρ_X is an isotropic, log-concave correlation function, then $\rho_{Y|X}^2 \leq \rho_{Y|\tilde{X}}^2$.

Proof: For convenience, write $\text{Var}(Y(s)) = \kappa^2$. Using the variance–covariance matrix of $(X(s), \tilde{X}(s), Y(s))^T$ given in (2.6), we then have

$$\begin{aligned}
 \rho_{Y|X}^2 &= \frac{\sigma_X^2 \beta_1^2}{\kappa^2} \rho_{X,\tilde{X}}^2(s, s) \\
 &= \frac{\sigma_X^2 \beta_1^2}{\kappa^2 K(s | \xi)^2} \int_{\mathcal{D}} \int_{\mathcal{D}} K(s, \mathbf{u} | \xi) \rho_X(\|\mathbf{s} - \mathbf{u}\|) K(s, \mathbf{v} | \xi) \rho_X(\|\mathbf{v} - \mathbf{s}\|) d\mathbf{u} d\mathbf{v} \\
 &\leq \frac{\sigma_X^2 \beta_1^2}{\kappa^2 K(s | \xi)^2} \int_{\mathcal{D}} \int_{\mathcal{D}} K(s, \mathbf{u} | \xi) K(s, \mathbf{v} | \xi) \rho_X(\|\mathbf{v} - \mathbf{u}\|) d\mathbf{u} d\mathbf{v} \\
 &= \frac{\sigma_X^2 \beta_1^2}{\kappa^2} \rho_{\tilde{X}}(s, s) \\
 &= \rho_{Y|\tilde{X}}^2,
 \end{aligned}$$

where the inequality comes from the log-concavity assumption. \square

The class of log-concave covariance functions includes the powered exponential, $\rho(s, s') = \exp\{-\phi\|\mathbf{s} - \mathbf{s}'\|^\alpha\}$ $0 \leq \alpha \leq 2$, by direct calculation. Also included are closed-form Matérn models, i.e., those with smoothness parameter ν of the form $\nu = k + 1/2$ for $k \in \{0, 1, 2, \dots\}$ again by direct calculation with an indication that this is the case for arbitrary ν (see Majumdar and Gelfand 2007). Also, it is easy to argue that convolution of covariance functions produces valid covariance functions. In fact, convolution of log-concave functions produces log-concave functions (for further examples, see, e.g., Gardner 2002).

Thus, for such covariance functions, using $X(s)$ instead of $\tilde{X}(s)$ results in less variation in $Y(s)$ explained. Additionally, notice that marginalizing over $\tilde{X}(s)$ adds extra spatial variation to $Y(s)$. To see this, notice that the covariance given by (2.9) has the added variance term $\beta_1^2 \sigma_X^2 \rho_{\tilde{X}}(s)$ (see the remark below expression (2.6) in this regard).

We now turn to Bayesian kriging at an unobserved location s_0 under the model given by (2.5). Bayesian kriging obtains the predictive distribution $[Y(s_0) | Y, X]$, where $Y = (Y(s_{y,1}), \dots, Y(s_{y,n}))^T$ is the vector of responses observed at spatial locations $s_{y,1}, \dots, s_{y,n}$, and $X = (X(s_{x,1}), \dots, X(s_{x,k}))^T$ is the vector of covariates observed at the locations $s_{x,1}, \dots, s_{x,k}$, where $\{s_{y,i}\}$ and $\{s_{x,j}\}$ are, potentially, misaligned (see Section 3). Mathematically,

$$\begin{aligned} [Y(s_0) | Y, X] &= \int_{\Theta} [Y(s_0) | \theta, Y, X] [\theta | Y, X] d\theta \\ &= \int_{\Theta} \int_{\mathbb{R}} [Y(s_0) | \tilde{X}(s_0), \theta, Y] [\tilde{X}(s_0) | \theta, X] \cdots \\ &\quad \times [\theta | Y, X] d\tilde{X}(s_0) d\theta, \end{aligned} \quad (2.10)$$

where $\theta \in \Theta$ are the model parameters. Based on (2.10), kriging under model (2.5) must occur in two stages. First, using $\theta \sim [\theta | Y, X]$, a kriged value for $\tilde{X}(s_0)$ is drawn from the conditional distribution $[\tilde{X}(s_0) | \theta, X]$. Due to the normality assumption on $X(s)$, $[\tilde{X}(s_0) | \theta, X]$ is also normal with mean and variance obtained using multivariate normal theory. Second, the kriged value of $\tilde{X}(s_0)$ is used to draw $Y(s_0) \sim [Y(s_0) | \tilde{X}(s_0), \theta, Y]$, where $[Y(s_0) | \tilde{X}(s_0), \theta, Y]$ is also Gaussian. Repeating the above steps M times will produce M draws $\{Y^{(1)}(s_0), \dots, Y^{(M)}(s_0)\}$ from $[Y(s_0) | Y, X]$.

As noted in the Introduction, we emphasize that the models proposed in this article are motivated by the objective of describing a local regression for $\mathbb{E}(Y(s))$ using neighboring covariate values. However, kriging under the model with kernel averaged predictors is not expected to differ much from that under the customary bivariate process spatial regression model. To see why this is the case, notice from (2.6) that the induced joint model for $(Y(s), X(s))$ is a GP with a mean that is linear in $X(s)$ and a covariance function which differs from that ordinarily used for spatial regression problems. Thus, co-kriging under the two models will only differ in the specification of the cross-covariance function. In practice, predictions under the two models will be similar as is demonstrated in Section 4.1 by Table 3.

On a related note, it is evident that model comparison between the local regression given $X(s)$ and given $\tilde{X}(s)$ cannot be made in predictive space; the models will be essentially indistinguishable. Rather, comparison resides in performance with regard to parametric inference using criteria such as mean square error and likelihood-based goodness-of-fit criterion such as AIC or DIC (see Section 4.1 for examples). In this regard, suppose that the kernel allows arbitrarily rapid decay, i.e., given s and s' such that $\|s - s'\| \leq \epsilon$, there exists a δ such that $K(s, s' | \xi) < \delta$. Then, we can view using $X(s)$ vs. $\tilde{X}(s)$ in the conditional mean for $Y(s)$ as reduced and full models, respectively, where the difference in dimension between the models is the dimension of ξ . But then, using within sample goodness of fit, the full model will always fit at least as well as the reduced model.

3. MODEL FITTING

Let $Y = (Y(s_{Y,1}), \dots, Y(s_{Y,n}))^T$ and $X = (X(s_{X,1}), \dots, X(s_{X,k}))^T$ be the vectors of observed responses and covariates, respectively. The locations of the observed responses $\{s_{y,i}\}_{i=1}^n$ need not be aligned with the locations of the observed covariates $\{s_{x,j}\}_{j=1}^k$ (e.g. ozone monitoring stations are much more sparse than, say, temperature monitoring stations). Thus, in the general case, the data are the two vectors Y and X with some locations in common as opposed to n $(Y(s_i), X(s_i))$ pairs.

Assuming that $\mu_X(s) = \mu_X$ for simplicity, a general Gibbs sampler for fitting (2.5) would proceed as follows:

1. sample $\tilde{X} = (\tilde{X}(s_1), \dots, \tilde{X}(s_n))^T \sim [\tilde{X} \mid -, X, Y]$,
2. sample $\mu_X \sim [\mu_X \mid -, \tilde{X}, X, Y]$,
3. sample $(\beta_0, \beta_1) \sim [\beta_0, \beta_1 \mid -, \tilde{X}, X, Y]$,
4. sample $(\xi, \phi_X, \sigma_X^2) \sim [\xi, \phi_X, \sigma_X^2 \mid -, \tilde{X}, X, Y]$,
5. sample $(\phi_Y, \sigma_Y^2, \tau_Y^2) \sim [\phi_Y, \sigma_Y^2, \tau_Y^2 \mid -, \tilde{X}, X, Y]$,

where “ $-$ ” denotes all other model parameters, and $\tilde{X} = (\tilde{X}(s_{y,1}), \dots, \tilde{X}(s_{y,n}))^T$ is the vector of $\{\tilde{X}(s_{y,i})\}_{i=1}^n$ aligned with the n locations where $Y(s)$ was observed. Steps 1, 2, and 3 above can be done directly as the complete conditional distributions for \tilde{X} , μ_X , and $(\beta_0, \beta_1)^T$ are all Gaussian, but steps 4 and 5 possibly require Metropolis–Hastings (MH) updates. For misaligned data, step 1 still has closed form because $(X, \tilde{X})^T$ follows a normal distribution with mean $\mathbf{1}_{k+n}\mu_X$ and covariance matrix

$$\text{Var} \begin{pmatrix} X \\ \tilde{X} \end{pmatrix} = \Sigma_{X\tilde{X}} = \sigma_X^2 \begin{pmatrix} \mathbf{R}_0 & \mathbf{R}_1 \\ \mathbf{R}_1^T & \mathbf{R}_2 \end{pmatrix}, \quad (3.1)$$

where $\mathbf{1}_n$ is a length n vector of ones, $\mathbf{R}_0 = \{\rho_X(s_{x,i}, s_{x,j})\}_{i,j=1}^k$, $\mathbf{R}_1 = \{\rho_{X,\tilde{X}}(s_{x,i}, s_{y,j})\}_{i=1, j=1}^{k,n}$, and $\mathbf{R}_2 = \{\rho_{\tilde{X}}(s_{y,i}, s_{y,j})\}_{i,j=1}^n$. Thus, because $(X, \tilde{X})^T$ is normally distributed, $[\tilde{X} \mid -, X]$ and $[\tilde{X} \mid -, X, Y]$ are also Gaussian distributions.

Computationally, step 1 above can be expensive for large n and/or k as it involves the inversion of an $n \times n$ and a $k \times k$ matrix which take order $\mathcal{O}(n^3)$ and $\mathcal{O}(k^3)$ operations. For discussion on handling such “large n ” problems see Banerjee et al. (2008).

Oftentimes, the values of $\tilde{X}(s)$ are not of direct interest. In these cases, $\tilde{X}(s)$ can be integrated out, and the user can fit the marginal model of $(X(s), Y(s))^T$. By working with the marginal model, the dimensionality of the posterior distribution is decreased by n . However, this dimension reduction comes at a price in that the complete conditional distribution for (β_0, β_1) is no longer available in closed form. Assuming that $\sigma_Y^2 = 0$ would also simplify computations because τ_Y^2 would then have a closed-form complete conditional distribution. If σ_Y^2 is assumed to be zero, the implication is that the marginal spatial correlation in $Y(s)$ is a result of spatial correlation in $X(s)$.

In general, the stochastic integrations in (2.6) cannot be done analytically, and, yet, evaluation of these integrals is required to either evaluate the likelihood or obtain a

draw from a complete conditional distribution. As suggested in Banerjee, Carlin, and Gelfand (2004), Monte Carlo (MC) integration becomes a good candidate for approximate computation of these integrals due to the relation $\rho_{X, \tilde{X}}(s, s') = \mathbb{E}_{\mathbf{u}}(\rho_X(s, \mathbf{u}))$, where $\mathbf{u} \sim K(s', \mathbf{u} | \xi) / K(s' | \xi)$. For kernels where \mathbf{u} can be sampled from easily, such as those in Table 1, Monte Carlo integration is computationally feasible but challenging for a few reasons. First, even for the easiest case of an aligned data set of n $(Y(s_i), X(s_i))$ pairs, a continuous prior for (ξ, ϕ_X) would require at least $n(n+1)$ such MC integrations at each iteration of the MCMC algorithm! This problem is compounded for cases where $\|s - s'\|$ is small and thousands, if not millions, of draws of $\mathbf{u} \sim K(s', \mathbf{u} | \xi) / K(s' | \xi)$ are required such that $\widehat{\text{Var}}(\tilde{X} | X)$ is positive definite. The remedy we adopt is to use a discrete prior for (ξ, ϕ_X) such that the associated variance–covariance matrices can be computed prior to the MCMC algorithm, saved, and the appropriate one retrieved at each iteration. Beyond computational convenience, specifying a discrete prior for (ξ, ϕ_X) is also theoretically warranted. Results in Zhang (2004) clarify that, with a weak prior on σ_X^2 , an informative prior is needed to identify ϕ_X .

Second, for complex regions \mathcal{D} , drawing $\mathbf{u} \sim K(s', \mathbf{u} | \xi) / K(s' | \xi)$ may be slow when the latter is a nonstandard distribution or a standard distribution truncated to \mathcal{D} . Additionally, $K(s | \xi)$ varies spatially, so a different distribution is sampled for each s_i for each iteration. Here, we take $\mathcal{D} = \mathbb{R}^2$ so that using disks, $K(s | \xi) = \pi r^2$ with $r = \xi g(\phi_X)$ for all $s \in \mathcal{D}$ and \mathbf{u} can be drawn from a uniform distribution on the disk with center s and radius r .

For well-behaved model fitting, K cannot have arbitrarily high concentration around s . Equivalently, choosing between models which use $X(s)$ or $\tilde{X}(s)$ in $\mathbb{E}(Y(s) | \theta, \tilde{X}(s), X(s))$ cannot be done with priors for ξ that allow K to have arbitrarily high concentration. To illustrate, let K be as defined in Section 2.1 with $\xi \in (0, 1)$. As ξ approaches zero, $\Sigma_{X\tilde{X}}$ (as in (3.1)) will be approximately singular leading to an unstable likelihood for $[X, \tilde{X} | \cdot]$. For the K we use, simulation experiments have shown that a lower bound of $\xi = 0.05$ works well in practice but more experience is needed to understand the sensitivity of the posterior to the prior on ξ for other choices of kernels.

4. APPLICATIONS OF THE SINGLE KERNEL AVERAGED PREDICTOR MODEL

This section applies the modeling of Section 2 to simulated and field data. The field data set consists of measurements of dew duration and shrub density over the Negev Desert in Israel and was previously analyzed by Banerjee and Gelfand (2002).

4.1. MODEL PERFORMANCE ON SIMULATED DATA

Data were simulated under model (2.5) with $K(s, s' | \xi) = \mathbb{I}_{\{\|s - s'\| \leq r\}}$ and $r = \xi g(\phi_X)$ as defined above, in a full factorial design for $n \in \{50, 100\}$, $\phi_X = 10$, $\xi \in \{0, 0.1, 0.3, 0.5\}$, $\sigma_X^2 = 1$, and $\beta = (0, 1)^T$. Values for τ_Y^2 were chosen such that the population coefficient of determination (ρ^2) was 0.3, 0.6, and 0.9, respectively, for the true model. We assumed that the locations of observed responses and predictors were aligned and $\{s_i, i = 1, \dots, n\}$

were confined to the unit square in \mathbb{R}^2 but \mathcal{D} was taken to be all of \mathbb{R}^2 so as to avoid the difficulty in dealing with locations near the boundary. The case of $\xi = 0$ implies $\mathbb{E}(Y(s) \mid \boldsymbol{\theta}, \tilde{X}, X) = \beta_0 + \beta_1 X(s)$. For convenience, σ_Y^2 was assumed to be zero. For each combination of (n, ξ, τ_Y^2) , 30 data sets were simulated with additional n values of $Y(s)$ left as a hold-out sample to determine predictive performance of the fitted models. Define $B(\hat{\beta}_1)$, $MSE(\hat{\beta}_1)$, $CIC(\beta_1)$, and $CIW(\beta_1)$ as the observed bias of the posterior mean $\hat{\beta}_1$, mean square error of $\hat{\beta}_1$, empirical 95% credible interval coverage for β_1 , and empirical 95% credible interval width for β_1 , respectively. Additionally, define $MSE(\hat{y}) = \sum_{j=1}^{30} \sum_{i=1}^n (y^{(j)}(s_i) - \hat{y}^{(j)}(s_i))^2 / (n \times 30)$ to be the average mean square error for the n training sample points, where $y^{(j)}(s_i)$ is the i th observation of the j th data set, and $\hat{y}^{(j)}(s_i)$ is the posterior mean of $\mathbb{E}(y^{(j)}(s_i) \mid \boldsymbol{\theta}, \tilde{X}, X)$. Finally, to assess predictive performance, define PMSE, PIC, and PIW as the average predictive mean square error, predictive interval coverage, and predictive interval width across all of the n hold-out values.

Two models were fit for comparison. The first model was the kernel averaged predictor (KAP) model given by (2.5), where $\sigma_Y^2 = 0$, and K was defined in Section 2.1. The second model was a point predictor (PtP) model with $\mathbb{E}(Y(s) \mid \boldsymbol{\theta}, \tilde{X}(s), X(s)) = \beta_0 + \beta_1 X(s)$, where $X(s)$ was defined in (2.1). The exponential correlation function with $\rho_X(s, s' \mid \phi_X) = \exp(-\phi_X \|s - s'\|)$ was used for both models. Discrete prior distributions for ϕ_X and ξ were used with mass at (30, 60, 300) and (0.1, 0.3, 0.5), respectively. Vague, but proper, conjugate prior distributions were assumed for the remaining parameters. Chains were run for an initial burn in period of 10,000 draws, and the following 10,000 were retained as draws from the posterior distribution.

To avoid redundancy, only the simulation results for the case of $\rho^2 = 0.9$ are displayed. Table 2 displays the results of the model fit along with values of DIC (Spiegelhalter et al. 2002) and the associated effective number of parameters (p_d). First, when the “true” model is $\mathbb{E}(Y(s) \mid \boldsymbol{\theta}, X, \tilde{X}) = \beta_0 + \beta_1 X(s)$, the KAP model has poor coverage ($CIC(\beta_1)$) of 33% and 3% for $n = 50$ and $n = 100$, respectively. However, the KAP model still has favorable values for $MSE(\hat{y})$. When $\xi > 0$, the KAP model performs well, while the performance of the PtP model declines rapidly as ξ increases. For example, notice that for the moderate value of $\xi = 0.1$, the coverage probability for the PtP model falls to 17% for $n = 50$ and 0% for $n = 100$. In terms of model choice, DIC correctly chooses the PtP model when $\xi = 0$ and the KAP model when $\xi > 0$. The simulation results for $\rho^2 \in \{0.3, 0.6\}$ were similar except that the discrepancy in DIC between the KAP and PtP models was smaller.

Table 3 displays the predictive results for the two models on the n hold-out values. The predictive performance of the models is similar in terms of PMSE and PIC. In terms of PIW, however, the KAP model seems to give smaller intervals. The comparable predictive performance between the two models was expected and clarified in Section 2.2.

4.2. KERNEL SELECTION SIMULATION STUDY

The simulation study here focuses on the problem of kernel selection. 50 data sets of $Y = (Y(s_{Y,1}), \dots, Y(s_{Y,n}))^T$ and $X = (X(s_{X,1}), \dots, X(s_{X,k}))^T$ with $n = 50$ and $k = 150$ (to introduce misalignment) were simulated from the KAP model (2.5) using $K(s, s' \mid$

Table 2. Model fit diagnostics of the kernel averaged predictor (KAP) and the point predictor (PtP) models on simulated data with $\rho^2 = 0.9$. The KAP model estimates β_1 well when $\xi > 0$.

n ξ		50				100			
		0	0.1	0.3	0.5	0	0.1	0.3	0.5
$B(\hat{\beta}_1)$	KAP	0.17	-0.02	-0.18	-0.23	0.15	-0.03	-0.16	-0.18
	PtP	0.01	-0.19	-0.43	-0.60	-0.01	-0.19	-0.46	-0.62
$MSE(\hat{\beta}_1)$	KAP	0.03	0.01	0.04	0.06	0.03	0.00	0.04	0.05
	PtP	0.00	0.04	0.19	0.36	0.00	0.04	0.21	0.38
$CIC(\beta_1)$	KAP	0.33	0.97	0.93	0.90	0.03	0.90	0.93	0.97
	PtP	0.97	0.17	0.00	0.00	0.93	0.00	0.00	0.00
$CIW(\beta_1)$	KAP	0.33	0.40	0.66	0.73	0.21	0.22	0.60	0.69
	PtP	0.22	0.27	0.28	0.26	0.15	0.17	0.19	0.18
$MSE(\hat{y})$	KAP	0.03	0.08	0.09	0.06	0.02	0.06	0.08	0.05
	PtP	0.11	0.19	0.21	0.20	0.11	0.19	0.22	0.21
p_d	KAP	22.65	18.95	17.25	19.34	52.87	43.96	39.84	46.32
	PtP	2.41	2.70	2.81	2.86	2.68	2.84	2.91	2.94
DIC	KAP	49.75	62.57	62.98	50.69	80.11	112.09	117.31	95.06
	PtP	37.07	64.57	69.95	65.59	68.51	121.66	138.35	130.28

Table 3. Predictive diagnostics of the kernel averaged predictor (KAP) and the point predictor (PtP) models on simulated data with $\rho^2 = 0.9$. The predictive properties of the KAP model are similar to those of the PtP model, but the KAP model gives slightly smaller predictive intervals.

n ξ		50				100			
		0	0.1	0.3	0.5	0	0.1	0.3	0.5
PMSE	KAP	1.12	0.86	0.56	0.38	1.08	0.91	0.54	0.40
	PtP	1.12	0.85	0.55	0.36	1.09	0.90	0.52	0.36
PIC	KAP	0.93	0.94	0.96	0.95	0.92	0.95	0.95	0.95
	PtP	0.95	0.95	0.97	0.96	0.96	0.96	0.96	0.95
PIW	KAP	3.89	3.61	2.91	2.39	3.69	3.51	2.74	2.26
	PtP	4.35	3.80	3.00	2.45	4.10	3.69	2.84	2.37

$\xi) = \mathbb{I}_{\{\|s-s'\| \leq r\}}$, where $r = \xi g(\phi_X)$, ρ_X is the exponential correlation function, $g(\phi_X) = 3/\phi_X$, $\phi_X = 10$, $\xi = 0.5$, $\sigma_X^2 = 1$, and $\beta = (0, 1)^T$. The value for τ_Y^2 was chosen such that $\rho_{Y|\tilde{X}}^2 = 0.5$. The locations for the $n = 50$ data points $\{s_{Y,1}, \dots, s_{Y,50}\}$ were drawn uniformly on the unit square in \mathbb{R}^2 . The $\{s_{X,j}\}$ were drawn uniformly on the larger square $(-0.5, 1.5) \times (-0.5, 1.5)$. Misalignment between the response and predictors is common in environmental studies but handled easily using the methods described in Section 3.

Three KAP models with kernels,

$$K_1(s, s' | \xi) = \mathbb{I}_{\{\|s-s'\| \leq \xi g(\phi_X)\}},$$

$$K_2(s, s' | \xi) = \frac{\phi_X}{\xi} \exp \left\{ -\frac{(s' - s)^T (s' - s)}{2\xi^2/\phi_X^2} \right\},$$

Table 4. Model fit diagnostics for the kernel selection simulation study. The kernels K_1 and K_2 produce similar results because they are of similar center and scale, but the misspecified kernel K_3 produces unsatisfactory results in terms of bias and credible interval coverage. This difference in performance is captured by both $\mathbb{E}(\rho_{Y|\tilde{X}}^2 | Y)$ and DIC.

Kernel	$B(\hat{\beta}_1)$	$CIC(\beta_1)$	$CIW(\beta_1)$	$\mathbb{E}(\rho_{Y \tilde{X}}^2 Y)$	DIC
$K_1(s, s' \xi)$	−0.12	0.96	0.84	0.46	91.95
$K_2(s, s' \xi)$	−0.21	0.92	0.90	0.46	93.82
$K_3(s, s' \xi)$	−0.63	0.50	1.06	0.29	103.72

and

$$K_3(s, s' | \xi) = \frac{\phi_X}{\xi} \exp \left\{ - \frac{(s' - (s + \theta))^T (s' - (s + \theta))}{2\xi^2/\phi_X^2} \right\},$$

were fit to each simulated data set, where $\xi \in (0, 1)$ for all kernels, and $\theta = (-0.2, 0.2)^T$. Notice that kernel K_1 is the true kernel and kernels K_2 and K_3 are simply Gaussian kernels with means s and $s + \theta$, respectively, and common standard deviation ξ/ϕ_X . The standard deviation of ξ/ϕ_X was chosen to confine the bulk of the mass of K_2 and K_3 to be within the spatial range $g(\phi_X)$ so as to prevent strong collinearity in the $\tilde{X}(s)$ (see Section 2.1). Notice that under this parameterization, the spatial range $g(\phi_X)$ is assumed to be, at most, 3 standard deviations away from the mean of the kernel in any direction. Each model assumed discrete priors for $[\phi_X]$ and $[\xi]$ with equal mass at (30, 10, 5) and (0.1, 0.5, 0.8), respectively. Here, θ is assumed to be known because we wish to investigate the consequences of using an inappropriate kernel. Table 4 displays $B(\hat{\beta}_1)$, $CIC(\beta_1)$, $CIW(\beta_1)$, $\mathbb{E}(\rho_{Y|\tilde{X}}^2 | Y)$, and DIC for each KAP model.

$\mathbb{E}(\rho_{Y|\tilde{X}}^2 | Y)$ and DIC can be used to select an appropriate kernel. That is, different kernels yield different predictors in the linear regression given by (2.5). Thus, comparing different kernels is essentially equivalent to comparing predictors. As such, DIC is a familiar criterion. As noted in Section 2, $\rho_{Y|\tilde{X}}^2 | Y$ is also an appropriate (but random) measure for the performance of a kernel averaged predictor model. Hence, its posterior mean provides a sensible kernel model comparison criterion.

The results in Table 4 give a few encouraging insights into kernel selection. First, similar kernels yield similar performance. Notice that KAP models using K_1 and K_2 exhibited acceptable performance despite the subtleties in the differences of their shapes while K_3 , which was shifted by a factor of θ , exhibited inferior bias, credible interval coverage, and credible interval width. Interestingly, when using K_2 , the value of ξ with the highest posterior probability was most frequently 0.8 across the data sets. Thus, despite the subtle difference in shape, the scale of K_2 was properly adjusted to capture the scale of the true kernel.

4.3. NEGEV DESERT FIELD DATA

The Negev Desert data set contains measurements of dew duration time in hundredths of an hour starting from 8 a.m. and shrub density (density of shrubs within a 5 meter \times

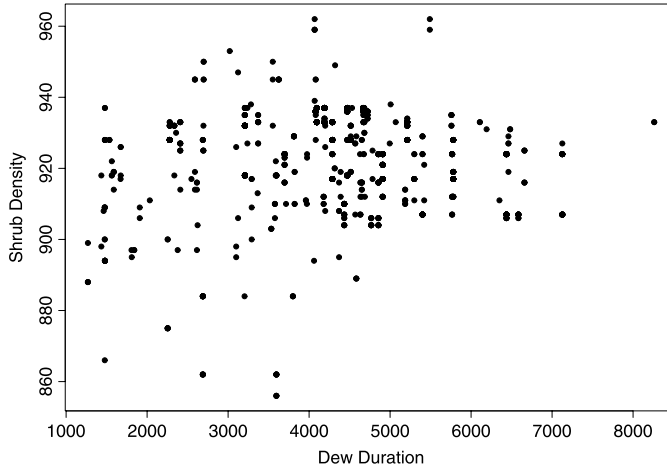


Figure 3. Scatter plot of dew duration by shrub density.

5 meter block) at $n = 1129$ aligned locations in the Negev Desert in Israel. As the Negev Desert is very arid, condensation can contribute considerably to the annual water levels. The analysis here is to determine the effect of shrub density on dew duration. Figure 3 displays the scatter plot of dew duration by shrub density, and Figure 4 displays contour plots of dew duration and shrub density over the region with sampling locations overlaid. All 1129 locations were used for the following analysis.

Let $Y(s_i)$ and $X(s_i)$ represent the measurement of dew duration and shrub density for UTM coordinate s_i , $i = 1, \dots, n$. Again, the KAP and PtP models described in the previous simulation study were fit to the shrub data. Because the primary interest of this study is the effect of shrub density on dew duration (β_1), both models assumed that the response process had no residual spatial correlation. A discrete uniform prior distribution for $g(\phi_X) = 3/\phi_X$ was used in both models, where mass was placed at 0.1, 0.3, and 0.5 times $\max_{(i,j)} \|s_i - s_j\|$. A discrete uniform prior distribution for ξ was used where masses were placed at the points (0.1, 0.2, 0.3, 0.4). Both models used vague, but proper, normal prior distributions for $[\beta_0, \beta_1]$ and vague, but proper, inverse-gamma distributions for $[\sigma_X^2, \tau_Y^2]$.

Markov chain Monte Carlo algorithms similar to that presented in Section 3 were run for an initial burn-in period of 10,000 iterations, and 20,000 draws were kept after this burn-in period. The KAP model took an average of 5.29 seconds per iteration compared to 1.66 seconds per iteration for the PtP model. Analysis of the chains using trace plots and the convergence diagnostics of Raftery and Lewis (1992) and Geweke (1992) showed that the chains had reasonable convergence properties and mixed well.

Comparing the KAP and PtP models, the KAP model had a DIC value of $\text{DIC} = 23,846$ versus $\text{DIC} = 25,111$ for the PtP model. The effective numbers of parameters were $p_d = 977.81$ and $p_d = 2.74$ for the KAP and PtP models, respectively.

Table 5 displays posterior quantiles for parameters that are common to both models, Table 6 lists the posterior distribution of $g(\phi_X)$ for both models, and the posterior distribution for ξ was $[\xi | X, Y] = (0.0, 0.97, 0.03, 0.0)$ for $\xi \in (0.1, 0.2, 0.3, 0.4)$, respectively. According to Table 5, all the quantiles of $[\beta_1 | X, Y]$ under the KAP model are less than

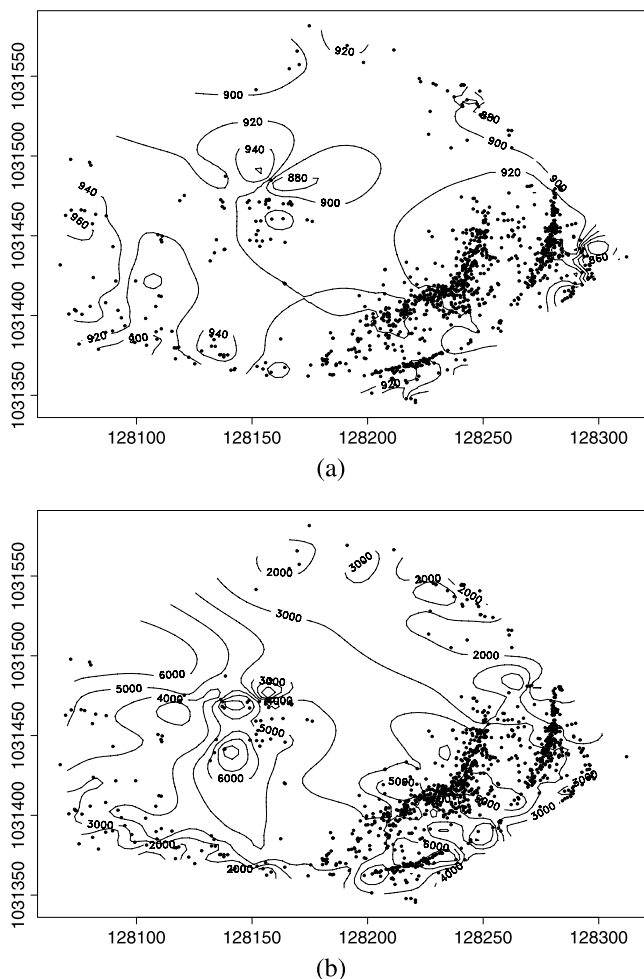


Figure 4. Contour plots of (a) dew duration and (b) shrub density over the Negev desert. The points correspond to the $n = 1129$ locations in universal transverse Mercator coordinates.

the associated quantiles under the PtP model. Thus, as expected, the effect of shrub density on dew duration is shrunk toward zero for the unblocked predictor model. Additionally, because the estimate of β_1 was shrunk toward zero, the estimate of β_0 differs between the two models as well. Notice also that $\rho_{\text{KAP}}^2 > \rho_{\text{PtP}}^2$, indicating that $\tilde{X}(s)$ is a better predictor of dew duration than $X(s)$.

For many ecological processes, response at a site would be expected to be influenced by covariate information from neighbors. The value of our approach is that we can incorporate this notion and quantify its benefit without actually observing the neighboring covariate levels. This is the primary finding of our reanalysis of the Negev desert data; that $\tilde{X}(s)$ was found to be a better predictor of $Y(s)$ than $X(s)$. This finding suggests new ecological interpretation with regard to the relationship between shrub density and dew duration within the Negev desert. Additionally, we find a stronger negative association between shrub density and dew duration than previous analyses had discovered. Though

Table 5. Posterior quantiles of parameters for the kernel averaged predictor (KAP) and point predictor (PtP) models used in the Negev Desert data application. As expected, $|\beta_1|$ and ρ^2 are greater for the KAP model.

Parameter	Model	Posterior Quantile				
		0.025	0.25	0.50	0.75	0.975
β_0	KAP	9.26e2	9.28e2	9.29e2	9.30e2	9.32e2
	PtP	9.25e2	9.26e2	9.27e2	9.28e2	9.29e2
β_1	KAP	-0.016	-0.012	-0.010	-0.008	-0.004
	PtP	-0.011	-0.008	-0.007	-0.005	-0.003
ρ^2	KAP	0.002	0.008	0.013	0.018	0.031
	PtP	0.002	0.006	0.009	0.014	0.024
μ_X	KAP	5.19e2	5.27e2	5.32e2	5.36e2	5.44e2
	PtP	5.19e2	5.27e2	5.32e2	5.36e2	5.44e2
σ_X^2	KAP	5.21e4	5.49e4	5.64e4	5.81e4	6.13e4
	PtP	5.28e4	5.56e4	5.73e4	5.89e4	6.23e4
τ_Y^2	KAP	2.59e2	2.73e2	2.81e2	2.89e2	3.05e2
	PtP	2.60e2	2.74e2	2.82e2	2.90e2	3.06e2

Table 6. Posterior distribution of the effective spatial range $g(\phi_X)$ for the shrub data using both a kernel averaged predictor (KAP) and point predictor (PtP) model.

Model	$g(\phi_X)$		
	24.82	74.44	124.07
KAP	0.93	0.04	0.03
PtP	0.33	0.34	0.33

this somewhat counterintuitive, the negative association between dew duration and shrub density is evident upon examination of the scatterplot in Figure 3.

5. EXTENSIONS

5.1. MULTIPLE KERNEL AVERAGED PREDICTORS

Extending the methodology of Section 2 to multiple predictors involves inducing a multivariate spatial process for $\mathbf{X}(s) = (X_1(s), \dots, X_p(s))^T$. Once a multivariate process has been specified for $\mathbf{X}(s)$, the extension follows by defining

$$\tilde{X}_j(s) = \frac{1}{K_j(s)} \int_{\mathcal{D}} K_j(s, \mathbf{u} \mid \boldsymbol{\xi}_j) X_j(\mathbf{u}) d\mathbf{u},$$

where $K_j(s, s' \mid \boldsymbol{\xi}_j)$ is the j th kernel with parameters $\boldsymbol{\xi}_j$ and extending (2.5) to be

$$\begin{aligned} Y(s) \mid \mathbf{X}(s), \tilde{\mathbf{X}}(s) &= \beta_0 + \sum_{j=1}^p \beta_j \tilde{X}_j(s) + \sigma_Y w_Y(s) + \epsilon_Y(s) \\ &= \beta_0 + \boldsymbol{\beta}^T \tilde{\mathbf{X}}(s) + \sigma_Y w_Y(s) + \epsilon_Y(s), \end{aligned} \quad (5.1)$$

where $w_Y(s)$ and $\epsilon_Y(s)$ are defined the same as in (2.5). It seems sensible that the choice of kernel would depend upon the predictor. For example, with K defined in Section 2.1, r should be allowed to vary with the covariate. Details and distributional results for the multiple kernel averaged predictor model are given in the supplementary materials where the multivariate spatial process for $X(s)$ is induced by the methods of coregionalization as described in Gelfand et al. (2004).

The multiple blocked predictor model in (5.1) poses significant computational challenges beyond those of the single kernel averaged predictor model in (2.5). Most noticeably, an MCMC algorithm to fit (5.1) requires that the length np vector $\tilde{X} = (\tilde{X}_1^T, \dots, \tilde{X}_p^T)^T$ be drawn at each iteration. Obtaining a draw of \tilde{X} involves inversion of a $kp \times kp$ and an $np \times np$ matrix which can be challenging even for moderately sized n , k , and p . Marginalization over $\tilde{X}(s)$ will remove the necessity to sample \tilde{X} at each iteration but conjugacy of full conditional distributions is lost. However, this marginalization allows a marked reduction in dimension of np . A problem which cannot be eluded via marginalization is the need to calculate integrals which are typically not available in closed form. Thus, to emphasize, a discrete prior for the parameters of the spatial correlation function and of the kernels $K_j(\cdot)$ becomes necessary for computational tractability.

We note that not all covariates are appropriate to kernel average. For instance, with categorical covariates, Gaussian process modeling will not be sensible, and kernel averaging will not provide an interpretable covariate level. Also, if trend surfaces are introduced to explain $Y(s)$, such surfaces would not be viewed as random; moreover, it may not be sensible to kernel average such surfaces.

5.2. GENERALIZED LINEAR MODELS WITH KERNEL AVERAGED PREDICTORS

Obviously, applications exist where the response variable is not continuous so that the use of a Gaussian distribution as a likelihood is inappropriate. The methods proposed herein naturally extend to noncontinuous $Y(s)$ through the use of link functions. Specifically, if $Y(s) \sim [Y(s) \mid \theta, \tilde{X}(s), X(s)]$, where $[\cdot]$ is some distribution function with parameters θ , then a generalized linear model with kernel averaged predictors would be $\mathbb{E}(Y(s) \mid \theta, \tilde{X}(s), X(s)) = h^{-1}(\beta_0 + \beta^T \tilde{X}(s))$, where h is an appropriate link function. For example, if $Y(s)$ is a count, then $[Y(s) \mid \theta, \tilde{X}(s), X(s)]$ is typically assumed to be the Poisson distribution, and $h^{-1}(\beta_0 + \beta^T \tilde{X}(s)) = \exp(\beta_0 + \beta^T \tilde{X}(s))$.

6. DISCUSSION AND FUTURE DIRECTIONS

We have presented an approach to account for neighboring predictor values in a spatial regression setting. The need for such an approach is apparent in that, to explain the response at a particular location, levels of certain predictors (e.g., precipitation, temperature) at neighboring locations can affect the response. The kernel averaged predictor (KAP) model elucidated in this article used kernel averaging to define a new predictor $\tilde{X}(s)$ for use in a spatial regression model. The parameters of the kernel control the extent to which neighboring predictors influence the response. An important result under the KAP model is that a regression model using $X(s)$ instead of $\tilde{X}(s)$ will diminish the effect of the predictor

on the response. This result was found to be the case for a real dataset of dew duration and shrub density over the Negev desert.

According to application, different kernels, perhaps oriented or location-shifted could be employed. Such kernels could be constructed to align with physical aspects such as meteorological conditions, elevation, or differential equations. For more complex kernels, more experience is needed to learn how well-associated parameters can be estimated.

Often, the MCMC algorithms used to estimate model parameters had to be run for a substantial amount of time in order to reach the limiting distribution. Analysis of the Markov chains run for the example above showed that an initial burn-in of 10,000 was necessary. For large n problems where each iteration requires the inversion of $n \times n$ matrices, the time needed to perform these iterations can compound quickly. Alternative computational methods such as reparameterization or, perhaps, Laplace approximations to the posterior as in Rue, Martino, and Chopin (2009) can be investigated.

Our development here has been exclusively for the static spatial setting. However, extension to incorporate dynamics—both explanatory and response variables recorded over time—opens a rich variety of process modeling opportunities. These possibilities are currently under investigation and will be reported on in future work.

SUPPLEMENTARY MATERIALS

Results for model misspecification when $\tilde{X}(s)$ is used when $X(s)$ is the appropriate predictor. (PDF)

Distributional results for the case of multiple predictors. (PDF)

ACKNOWLEDGEMENTS

This work was supported in part by NIH grant number 5R01-ES-014843-01A2. Any opinions, findings, and conclusions or recommendations expressed in this publications are those of the authors and do not necessarily reflect the views of the National Institute of Health. This research was also supported in part by NSF grant DMS-0914906 to the National Institute of Statistical Sciences. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the authors and do not necessarily reflect the views of the National Science Foundation.

[Published Online January 2011.]

REFERENCES

- Baíllo, A., and Grané, A. (2009), "Local Linear Regression for Functional Predictor and Scalar Response," *Journal of Multivariate Analysis*, 100, 102–111.
- Banerjee, S., and Gelfand, A. E. (2002), "Prediction, Interpolation, and Regression for Spatially Misaligned Data Points," *Sankhya A*, 64, 227–245.
- Banerjee, S., Carlin, B. P., and Gelfand, A. E. (2004), *Hierarchical Modeling and Analysis for Spatial Data*, Boca Raton: Chapman and Hall/CRC.
- Banerjee, S., Gelfand, A. E., Finley, A. O., and Sang, H. (2008), "Gaussian Predictive Process Models for Large Spatial Datasets," *Journal of the Royal Statistical Society, Series B*, 70, 825–848.

- Berliner, L. M. (2000), "Hierarchical Bayesian Modeling in the Environmental Sciences," *Journal of the German Statistical Society*, 84.
- Chiles, J. P., and Delfiner, P. (1999), *Geostatistics: Modeling Spatial Uncertainty*, New York: Wiley Interscience.
- Cressie, N. A. C. (1993), *Statistics for Spatial Data*, New York: Wiley.
- Diggle, P. J., Tawn, J. A., and Moyeed, R. A. (1998), "Model-Based Geostatistics," *Applied Statistics*, 47, 299–350.
- Gardner, R. J. (2002), "The Brunn–Minkowski Inequality," *Bulletin of the American Mathematical Society*, 39, 355–405.
- Gelfand, A. E., Kim, H. J., Sirmans, C. F., and Banerjee, S. (2003), "Spatial Modeling With Spatially Varying Coefficient Processes," *Journal of the American Statistical Association*, 98, 387–396.
- Gelfand, A. E., Schmidt, A. M., Banerjee, S., and Sirmans, C. F. (2004), "Nonstationary Multivariate Process Modeling Through Spatially Varying Coregionalization" (with discussion), *Test*, 12, 1–50.
- Geweke, J. (1992), "Evaluating the Accuracy of Sampling-Based Approaches to the Calculation of Posterior Moments," in *Bayesian Statistics 4*, eds. J. M. Bernardo, J. Berger, A. P. Dawid, and A. F. M. Smith, Oxford: Oxford University Press, pp. 169–194.
- Higdon, D. (1998), "A Process-Convolution Approach to Modelling Temperatures in the North Atlantic Ocean," *Environmental and Ecological Statistics*, 5, 173–190.
- Majumdar, A., and Gelfand, A. E. (2007), "Multivariate Spatial Modeling for Geostatistical Data Using Convolved Covariance Functions," *Journal of Mathematical Geology*, 39, 225–245.
- Matérn, B. (1986), *Spatial Variation* (2nd ed.), Berlin: Springer.
- Raftery, A. E., and Lewis, S. (1992), "How Many Iterations in the Gibbs Sampler," in *Bayesian Statistics 4*, eds. J. M. Bernardo, J. Berger, A. P. Dawid, and A. F. M. Smith, Oxford: Oxford University Press, pp. 763–773.
- Royle, J. A., and Berliner, L. M. (1999), "A Hierarchical Approach to Multivariate Spatial Modeling and Prediction," *Journal of Agricultural, Biological, and Environmental Statistics*, 4.
- Rue, H., Martino, S., and Chopin, N. (2009), "Approximate Bayesian Inference for Latent Gaussian Models by Using Integrated Nested Laplace Approximations" (with discussion), *Journal of the Royal Statistical Society, Series B*, 71, 319–392.
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P., and van der Linde, A. (2002), "Bayesian Measures of Model Complexity and Fit" (with discussion), *Journal of the Royal Statistical Society, Series B*, 64, 583–639.
- Wackernagel, H. (2003), *Multivariate Geostatistics*, Berlin: Springer.
- Wikle, C. K. (2002), "A Kernel Based Spectral Method for Non-Gaussian Spatio-Temporal Processes," *Statistical Modelling*, 2, 299–314.
- Xu, K., Wikle, C. K., and Fox, N. L. (2005), "A Kernel-Based Spatio-Temporal Dynamical Model for Nowcasting Weather Radar Reflectivities," *Journal of the American Statistical Association*, 100, 1133–1144.
- Zhang, H. (2004), "Inconsistent Estimation and Asymptotically Equal Interpolations in Model-Based Geostatistics," *Journal of the American Statistical Association*, 99, 250–261.