

spmodel: spatial statistical modeling and prediction in R

Michael Dumelle ¹*, Matt Higham ², Jay M. Ver Hoef ³

¹ United States Environmental Protection Agency, 200 SW 35th St, Corvallis, OR, 97333

² St. Lawrence University Department of Math, Computer Science, and Statistics

³ National Oceanic and Atmospheric Administration Alaska Fisheries Science Center, Marine Mammal Laboratory, Seattle, WA, 98115

* Corresponding author: Dumelle.Michael@epa.gov

Abstract

`spmodel` is an **R** package used to fit, summarize, and predict for a variety spatial statistical models. Parameters are estimated using various methods. Additional modeling features include anisotropy, random effects, partition factors, big data approaches, and more. Model-fit statistics are used to summarize, visualize, and compare models. Predictions at unobserved locations are readily obtainable. This manuscript corresponds to `spmodel` version 0.1.0.

Introduction

Spatial data are ubiquitous in everyday life and the scientific literature. As such, it is becoming increasingly important to properly analyze spatial data. Spatial data can be analyzed using a statistical model that explicitly incorporates the spatial dependence among observations. Incorporating this spatial dependence can be challenging, but ignoring it often yields poor statistical models that incorrectly quantify uncertainty, which impacts the validity of hypothesis tests, confidence intervals, and predictions intervals. `spmodel` provides tools to easily incorporate spatial dependence, building upon commonly used **R** functions like `lm()`.

`spmodel` implements a model-based inference, which relies on fitting a statistical model. This is a different way to analyze data than design-based inference, which relies on random sampling and estimators that incorporate the properties of the random sample. Spatial random sampling and design-based inference can be performed using the `spsurvey` **R** package [1]. Though not explicitly required, random sampling can still benefit model-based inference when analyzing spatial data [2].

[3] defines two types of spatial data that can be analyzed from a model-based perspective: point-referenced data and areal data (sometimes called lattice data). Spatial data are point-referenced when they are observed at point-locations indexed by x-coordinates and y-coordinates on a spatially continuous surface with an infinite number of locations. Spatial models fit to point-referenced data are sometimes called geostatistical models. Spatial data are areal when they are observed as part of a finite network of polygons whose connections are indexed by a neighborhood structure. For example, the polygons may represent counties in a state who are neighbors if they share at least one boundary. Spatial models fit to areal data are sometimes called spatial autoregressive models.

Several **R** packages for analyzing point-referenced and areal spatial data exist. For point-referenced data, they include `fields` [4], `FRK` [5], `geoR` [6], `GpGp` [7], `gstat` [8], `LatticeKrig` [9], `spatial` [10], `spBayes` [11], and `spNNGP` [12]. For areal data, they include `CARBayes` [13] and `hglm` [14]. Unlike these aforementioned packages, `spmodel` is designed to analyze both point-referenced and areal data using a common framework and syntax structure. `spmodel` also offers many features missing from the aforementioned **R** packages. Together in one **R** package, `spmodel` offers detailed model summaries, extensive model diagnostics, random effects, anisotropy, big data methods, prediction, and more.

The rest of this article is organized as follows. In Section , we give a brief theoretical introduction to spatial linear models. In Section , we outline the variety of methods used to estimate the parameters of spatial linear models. In Section , we explain how to obtain predictions at unobserved locations. In Section , we detail some advanced modeling features, including random effects, partition factors, anisotropy, and big data approaches. In Section , we end with a short discussion. Before proceeding, we load `spmodel` by running

```
R> library(spmodel)
```

We will create visualizations using `ggplot2` [15], which we load by running

```
R> library(ggplot2)
```

We will also show code that can be used to create interactive visualizations of spatial data with `mapview` [16]. `mapview` also has many backgrounds available that contextualize spatial data with topographical information. Before running the `mapview` code interactively, make sure `mapview` is loaded.

The Spatial Linear Model

Statistical linear models are often parameterized as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (1)$$

where for a sample size n , \mathbf{y} is an $n \times 1$ column vector of response variables, \mathbf{X} is an $n \times p$ design (model) matrix of explanatory variables, $\boldsymbol{\beta}$ is an $p \times 1$ column vector of fixed effects controlling the impact of \mathbf{X} on \mathbf{y} , and $\boldsymbol{\epsilon}$ is an $n \times 1$ column vector of random errors. We typically assume that $E(\boldsymbol{\epsilon}) = \mathbf{0}$ and $\text{Cov}(\boldsymbol{\epsilon}) = \sigma_{\epsilon}^2 \mathbf{I}$, where $E(\cdot)$ denotes expectation, $\text{Cov}(\cdot)$ denotes covariance, σ_{ϵ}^2 denotes a variance parameter, and \mathbf{I} denotes the identity matrix.

The model in Equation 1 assumes the elements of \mathbf{y} are uncorrelated. Typically for spatial data, elements of \mathbf{y} are correlated, as observations close together in space tend to be more similar than observations far apart [17]. Failing to properly accommodate the spatial dependence in \mathbf{y} can cause researchers to draw incorrect conclusions about their data. To accommodate spatial dependence in \mathbf{y} , an $n \times 1$ spatial random effect, $\boldsymbol{\tau}$, is added to Equation 1, yielding the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\tau} + \boldsymbol{\epsilon}, \quad (2)$$

where $\boldsymbol{\tau}$ is independent of $\boldsymbol{\epsilon}$, $E(\boldsymbol{\tau}) = \mathbf{0}$, $\text{Cov}(\boldsymbol{\tau}) = \sigma_{\tau}^2 \mathbf{R}$, and \mathbf{R} is a matrix that determines the spatial dependence structure in \mathbf{y} and depends on a range parameter, ϕ . We discuss \mathbf{R} in more detail shortly. The parameter σ_{τ}^2 is called the spatially dependent random error variance or partial sill. The parameter σ_{ϵ}^2 is called the spatially independent random error variance or nugget. These two variance parameters are

henceforth more intuitively written as σ_{de}^2 and σ_{ie}^2 , respectively. The covariance of \mathbf{y} is denoted $\mathbf{\Sigma}$ and given by $\sigma_{de}^2 \mathbf{R} + \sigma_{ie}^2 \mathbf{I}$. The parameters that compose this covariance are typically referenced by the vector $\boldsymbol{\theta}$, which is called the covariance parameter vector.

Equation 2 is called the spatial linear model. The spatial linear model applies to both point-referenced and areal data. The `splm()` function is used to fit spatial linear models for point-referenced data (these are often called geostatistical models). One spatial covariance function available in `splm()` is the exponential spatial covariance function, which has an \mathbf{R} matrix given by

$$\mathbf{R} = \exp(-\mathbf{M}/\phi),$$

where \mathbf{M} is a matrix of Euclidean distances among observations. Recall that ϕ is the range parameter, controlling the behavior of \mathbf{R} as a function of distance. Parameterizations for other `splm()` spatial covariance types and their \mathbf{R} matrices can be seen by running `help("splm", "spmodel")` or `vignette("technical", "spmodel")`. Some of these spatial covariance types (e.g., Matérn) depend on an extra parameter beyond σ_{de}^2 , σ_{ie}^2 , and ϕ .

The `spautorm()` function is used to fit spatial linear models for areal data (these are often called spatial autoregressive models). One spatial autoregressive covariance function available in `spautorm()` is the simultaneous autoregressive spatial covariance function, which has an \mathbf{R} matrix given by

$$\mathbf{R} = [(\mathbf{I} - \phi \mathbf{W})(\mathbf{I} - \phi \mathbf{W})^\top]^{-1},$$

where \mathbf{W} is a weight matrix describing the neighborhood structure in \mathbf{y} . Parameterizations for `spautorm()` spatial covariance types and their \mathbf{R} matrices can be seen by running `help("spautorm", "spmodel")` or `vignette("technical", "spmodel")`.

One way to define \mathbf{W} is through queen contiguity [18]. Two observations are queen contiguous if they share a boundary. The ij th element of \mathbf{W} is then one if observation i and observation j are queen contiguous and zero otherwise. Observations are not considered neighbors with themselves, so each diagonal element of \mathbf{W} is zero.

Sometimes each element in the weight matrix \mathbf{W} is divided by its respective row sum. This is called row-standardization. Row-standardizing \mathbf{W} has several benefits, which are discussed in detail by [19].

Model Fitting

In this section, we show how to use the `splm()` and `spautorm()` functions to estimate parameters of the spatial linear model. We also explore diagnostic tools in `spmodel` that evaluate model fit. The `splm()` and `spautorm()` functions share similar syntactic structure with the `lm()` function used to fit non-spatial linear models (linear models without spatial dependence) from Equation 1. `splm()` and `spautorm()` generally require at least three arguments:

- **formula**: a formula that describes the relationship between the response variable (\mathbf{y}) and explanatory variables (\mathbf{X})
 - **formula** in `splm()` is the same as **formula** in `lm()`
- **data**: a `data.frame` or `sf` object that contains the response variable, explanatory variables, and spatial information
- **spcov_type**: the spatial covariance type ("`exponential`", "`matern`", "`car`", etc)

If `data` is an `sf` [20] object, spatial information is stored in the object's geometry. If `data` is a `data.frame`, then the x-coordinates and y-coordinates must be provided via the `xcoord` and `ycoord` arguments (for point-referenced data) or the weight matrix must be provided via the `W` argument (for areal data).

In the following subsections, we use the point-referenced `moss` data, an `sf` object that contains data on heavy metals in mosses near a mining road in Alaska. We view the first few rows of `moss` by running

```
R> moss

Simple feature collection with 365 features and 7 fields
Geometry type: POINT
Dimension: XY
Bounding box: xmin: -445884.1 ymin: 1929616 xmax: -383656.8 ymax: 2061414
Projected CRS: NAD83 / Alaska Albers
# A tibble: 365 x 8
  sample field_dup lab_rep year sideroad log_dist2road log_Zn      geometry
  <fct>   <fct>     <fct> <fct> <fct>          <dbl> <dbl>   <POINT>
1 001PR 1         1      2001 N             2.68  7.33 (-413585.3 1997
2 001PR 1         2      2001 N             2.68  7.38 (-413585.3 1997
3 002PR 1         1      2001 N             2.54  7.58 (-415367.2 1996
4 003PR 1         1      2001 N             2.97  7.63 (-417186.1 1995
5 004PR 1         1      2001 N             2.72  7.26 (-411756.6 1999
6 005PR 1         1      2001 N             2.76  7.65 (-419008.8 1994
7 006PR 1         1      2001 S             2.30  7.59 (-420822.6 1994
8 007PR 1         1      2001 N             2.78  7.16 (-404383.2 2007
9 008PR 1         1      2001 N             2.93  7.19 (-409964.1 2003
10 009PR 1         1      2001 N             2.79  8.07 (-424469.7 1992
# ... with 355 more rows
```

We can learn more about `moss` by running `help("moss", "spmodel")`, and we can visualize the distribution of log zinc concentration in `moss` (Figure 1) by running

```
R> ggplot(moss, aes(color = log_Zn)) +
+   geom_sf(size = 2) +
+   scale_color_viridis_c() +
+   theme_gray(base_size = 14)
```

Estimation

Generally the covariance parameters (θ) and fixed effects (β) of the spatial linear model require estimation. The default estimation method in `spmodel` is restricted maximum likelihood [21–23]. Maximum likelihood estimation is also available. For point-referenced data, semivariogram weighted least squares [24] and semivariogram composite likelihood [25] are additional estimation methods. The estimation method is chosen using the `estmethod` argument.

We estimate parameters of a spatial linear model regressing log zinc concentration (`log_Zn`) on log distance to a haul road (`log_dist2road`) using an exponential spatial covariance function by running

```
R> spmod <- splm(log_Zn ~ log_dist2road, moss, spcov_type = "exponential")
```

We summarize the model fit by running

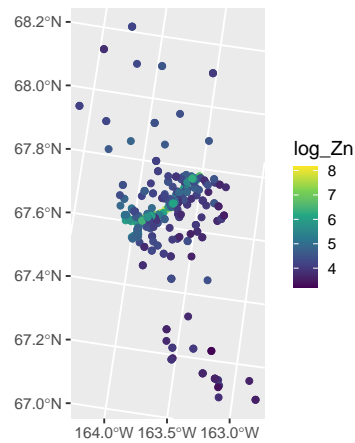


Fig 1. Distribution of log zinc concentration in the moss data.

```
R> summary(spmod)
146

Call:
147
splm(formula = log_Zn ~ log_dist2road, data = moss, spcov_type = "exponential")
149

Residuals:
150
    Min       1Q   Median       3Q      Max
151 -2.6801 -1.3606 -0.8103 -0.2485  1.1298
152
153
Coefficients (fixed):
154
              Estimate Std. Error z value Pr(>|z|)
155 (Intercept)   9.76825    0.25216   38.74  <2e-16 ***
156 log_dist2road -0.56287    0.02013  -27.96  <2e-16 ***
157 ---
158
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
159
160
Pseudo R-squared: 0.683
161
162
Coefficients ( spatial covariance):
163
           de           ie        range
164 3.595e-01 7.897e-02 8.237e+03
165
```

The fixed effects coefficient table contains estimates, standard errors, z-statistics, and asymptotic p-values for each fixed effect. From this table, we notice there is evidence that mean log zinc concentration significantly decreases with distance from the haul road (p-value < 2e-16). We see the fixed effect estimates by running

```
R> coef(spmod)
170

(Intercept) log_dist2road
171 9.7682525    -0.5628713
172
```

The model summary also contains the exponential spatial covariance parameter estimates, which we can view by running

```
R> coef(spmod, type = "spcov")
175
```

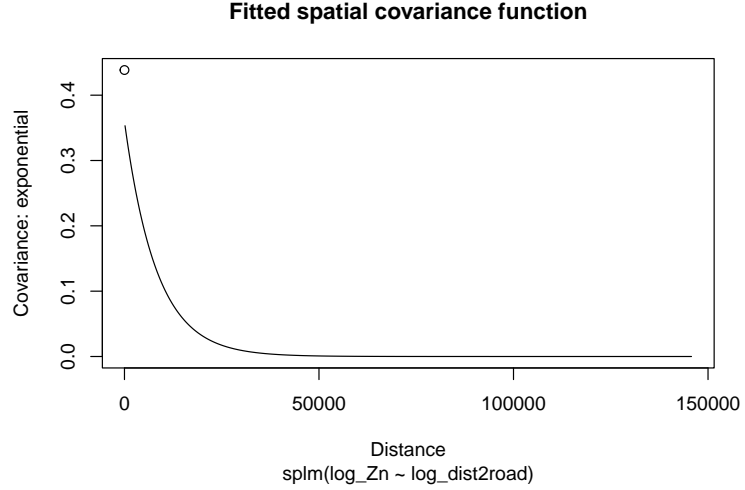


Fig 2. Empirical spatial covariance of fitted model. The open circle at a distance of zero represents the $\sigma_{de}^2 + \sigma_{ie}^2$. The solid line at positive distances represents the product of σ_{de}^2 and the correlation function at a particular distance.

```

              de          ie          range          rotate          scale
3.595316e-01 7.896824e-02 8.236712e+03 0.000000e+00 1.000000e+00
attr(,"class")
[1] "exponential"

```

The dependent random error variance (σ_{de}^2) is estimated to be approximately 0.36 and the independent random error variance (σ_{ie}^2) is estimated to be approximately 0.079. The range (ϕ) is estimated to be approximately 8,237. The effective range is the distance at which the spatial covariance is approximately zero. For the exponential covariance, the effective range is 3ϕ . This means that observations whose distance is greater than 24,711 meters are approximately uncorrelated. The **rotate** and **scale** parameters affect the modeling of anisotropy (Section). By default, they are assumed to be zero and one, respectively, which means that anisotropy is not modeled (i.e., the spatial covariance is assumed isotropic, or independent of direction). We plot the fitted spatial covariance function (Figure 2) by running

```
R> plot(spmo, which = 7)
```

We can learn more about the plots available for fitted models by running `help("plot.spmo", "spmodel")`.

Model-Fit Statistics

The quality of model fit can be assessed using a variety of statistics readily available in `spmodel`. The first model-fit statistic we consider is the pseudo R-squared. The pseudo R-squared is a generalization of the classical R-squared from non-spatial linear models that quantifies the proportion of variability in the data explained by the fixed effects. The pseudo R-squared is defined as

$$PR2 = 1 - \frac{\mathcal{D}(\hat{\Theta})}{\mathcal{D}(\hat{\Theta}_0)},$$

where $\mathcal{D}(\hat{\Theta})$ is the deviance of the fitted model indexed by parameter vector $\hat{\Theta}$ and $\mathcal{D}(\hat{\Theta}_0)$ is the deviance of an intercept-only model indexed by parameter vector $\hat{\Theta}_0$. For maximum likelihood, $\hat{\Theta} = \{\hat{\theta}, \hat{\beta}\}$. For restricted maximum likelihood $\hat{\Theta} = \{\hat{\theta}\}$.

We compute the pseudo R-squared by running

```
R> pseudoR2(spmod)
[1] 0.6829687
```

Roughly 68% of the variability in log zinc is explained by log distance from the road. The pseudo R-squared can be adjusted to account for the number of explanatory variables using the `adjust` argument. Pseudo R-squared (and the adjusted version) is most helpful for comparing models that have the same covariance structure.

The next two model-fit statistics we consider are the spatial AIC and AICc [26]. The AIC and AICc evaluate the fit of a model with a penalty for the number of parameters estimated. This penalty balances model fit and model parsimony. Lower AIC and AICc indicate a better balance of model fit and parsimony. The AICc is a correction to AIC for small sample sizes. As the sample size increases, AIC and AICc converge.

The spatial AIC and AICc are given by

$$\begin{aligned} \text{AIC} &= -2\ell(\hat{\Theta}) + 2(|\hat{\Theta}|) \\ \text{AICc} &= -2\ell(\hat{\Theta}) + 2n(|\hat{\Theta}|)/(n - |\hat{\Theta}| - 1), \end{aligned}$$

where $\ell(\hat{\Theta})$ is the log-likelihood of the data evaluated at the estimated parameter vector $\hat{\Theta}$ that maximized $\ell(\Theta)$, $|\hat{\Theta}|$ is the cardinality of $\hat{\Theta}$, and n is the sample size. As with the deviance, for maximum likelihood, $\hat{\Theta} = \{\hat{\theta}, \hat{\beta}\}$, and for restricted maximum likelihood $\hat{\Theta} = \{\hat{\theta}\}$. There are some nuances to consider when comparing AIC across models: AIC comparisons between a model fit using restricted maximum likelihood and a model fit using maximum likelihood are meaningless, as the models are fit with different likelihoods; AIC comparisons between models fit using restricted maximum likelihood are only valid when the models have the same fixed effect structure; AIC comparisons between models fit using maximum likelihood are valid even when the models have different fixed effect structures [27].

Suppose we want to quantify the difference in model quality between the spatial model and a non-spatial model using the AIC and AICc criteria. We fit a non-spatial model (Equation 1) in `spmodel` by running

```
R> lmod <- splm(log_Zn ~ log_dist2road, moss, spcov_type = "none")
```

We compute the spatial AIC and AICc of the spatial model and non-spatial model by running

```
R> AIC(spmod, lmod)
```

```
      df      AIC
spmod  3 373.2089
lmod   1 636.0635
```

```
R> AICc(spmod, lmod)
```

```
      df      AICc
spmod  3 373.2754
lmod   1 636.0745
```

The noticeably lower AIC and AICc of the spatial model indicate that it is a better fit to the data than the non-spatial model.

Another approach to comparing the fitted models is to perform leave-one-out cross validation [28]. In leave-one-out cross validation, a single observation is removed from the data, the model is re-fit, and a prediction is made for the held-out observation. Then, a loss metric like mean-squared-prediction error is computed and used to evaluate model fit. The lower the mean-squared-prediction error, the better the model fit. For computational efficiency, leave-one-out cross validation in `spmodel` is performed by first estimating θ using all the data and then re-estimating β for each observation. We perform leave-one-out cross validation for the spatial and non-spatial model by running

```
R> loocv(spmod)
```

```
[1] 0.1110895
```

```
R> loocv(lmod)
```

```
[1] 0.3237897
```

The noticeably lower mean-squared-prediction error of the spatial model indicates that it is a better fit to the data than the non-spatial model.

Diagnostics

In addition to model fit metrics, `spmodel` provides functions to compute diagnostic metrics that help assess model assumptions and identify unusual observations.

An observation is said to have high leverage if its combination of explanatory variable values is far from the mean vector of the explanatory variables. For a non-spatial model, the leverage of the i th observation is the i th diagonal element of the hat matrix given by

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top.$$

For a spatial model, the leverage of the i th observation is the i th diagonal element of the spatial hat matrix given by

$$\mathbf{H}^* = (\mathbf{X}^*(\mathbf{X}^{*\top} \mathbf{X})^{-1} \mathbf{X}^{*\top}),$$

where $\mathbf{X}^* = \Sigma^{-1/2} \mathbf{X}$ and $\Sigma^{-1/2}$ is the inverse matrix square root of the covariance matrix, Σ [29]. The spatial hat matrix can be viewed as the non-spatial hat matrix applied \mathbf{X}^* instead of \mathbf{X} . We compute the hat values (leverage) by running

```
R> hatvalues(spmod)
```

Larger hat values indicate more leverage.

The fitted value of an observation is the estimated mean response given the observation's explanatory variable values and the model fit:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}.$$

We compute the fitted values by running

```
R> fitted(spmod)
```


Fitted values for the spatially dependent random errors ($\boldsymbol{\tau}$), spatially independent random errors ($\boldsymbol{\epsilon}$), and random effects can also be obtained via `fitted()` by changing the `type` argument.

The residuals measure each response's deviation from its fitted value. The raw residuals are given by

$$\mathbf{e}_r = \mathbf{y} - \hat{\mathbf{y}}.$$

We compute the raw residuals of the spatial model by running

```
R> residuals(spmod)
```

The raw residuals are typically not directly checked for linear model assumptions, as they have covariance closely resembling the covariance of \mathbf{y} . Pre-multiplying the residuals by $\boldsymbol{\Sigma}^{-1/2}$ yields the Pearson residuals [30]:

$$\mathbf{e}_p = \boldsymbol{\Sigma}^{-1/2} \mathbf{e}_r.$$

When the model is correct, the Pearson residuals have mean zero, variance approximately one, and are uncorrelated. We compute the Pearson residuals of the spatial model by running

```
R> residuals(spmod, type = "pearson")
```

The covariance of \mathbf{e}_p is $(\mathbf{I} - \mathbf{H}^*)$, which is approximately \mathbf{I} for large sample sizes. Explicitly dividing \mathbf{e}_p by the respective diagonal element of $(\mathbf{I} - \mathbf{H}^*)$ yields the standardized residuals [30]:

$$\mathbf{e}_s = \frac{\mathbf{e}_p}{\sqrt{(1 - \text{diag}(\mathbf{H}^*))}},$$

where $\text{diag}(\mathbf{H}^*)$ denotes the diagonal of \mathbf{H}^* .

We compute the standardized residuals of the spatial model by running

```
R> residuals(spmod, type = "standardized")
```

or

```
R> rstandard(spmod)
```

When the model is correct, the standardized residuals have mean zero, variance one, and are uncorrelated. It is common to check linear model assumptions through visualizations. We can plot the standardized residuals vs fitted values by running

```
R> plot(spmod, which = 1) # figure omitted
```

When the model is correct, the standardized residuals should be evenly spread around zero with no discernible pattern. We can plot a normal QQ-plot of the standardized residuals by running

```
R> plot(spmod, which = 2) # figure omitted
```

When the standardized residuals are normally distributed, they should closely follow the normal QQ-line.

An observation is said to be influential if its omission has a large impact on model fit. Typically, this is measured using Cook's distance [31]. For the non-spatial model, the Cook's distance of the i th observation is denoted \mathbf{D} and given by

$$\mathbf{D} = \mathbf{e}_s^2 \frac{\text{diag}(\mathbf{H})}{p(1 - \text{diag}(\mathbf{H}))},$$

where p is the dimension of β (the number of fixed effects).

For a spatial model, the Cook's distance of the i th observation is denoted \mathbf{D}^* and given by

$$\mathbf{D}^* = \mathbf{e}_s^2 \frac{\text{diag}(\mathbf{H}^*)}{p(1 - \text{diag}(\mathbf{H}^*))}.$$

A larger Cook's distance indicates more influence from the observation. We compute Cook's distance by running

```
R> cooks.distance(spmod)
```

The Cook's distance versus leverage (hat values) can be visualized by running

```
R> plot(spmod, which = 6) # figure omitted
```

Though we described the model diagnostics in this subsection using Σ , generally the covariance parameters are estimated and Σ is replaced with $\hat{\Sigma}$.

The broom functions: tidy(), glance(), and augment()

The tidy(), glance(), and augment() functions from the broom R package [32] provide convenient output for many of the model fit and diagnostic metrics discussed in the previous two sections. The tidy() function returns a tidy tibble of the coefficient table from summary():

```
R> tidy(spmod)
```

```
# A tibble: 2 x 5
```

	term	estimate	std.error	statistic	p.value
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
1	(Intercept)	9.77	0.252	38.7	0
2	log_dist2road	-0.563	0.0201	-28.0	0

This tibble format makes it easy to pull out the coefficient names, estimates, standard errors, z-statistics, and p-values from the summary() output.

The glance() function returns a tidy tibble of model-fit statistics:

```
R> glance(spmod)
```

```
# A tibble: 1 x 9
```

	n	p	npar	value	AIC	AICc	logLik	deviance	pseudo.r.squared
	<int>	<dbl>	<int>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	365	2	3	367.	373.	373.	-184.	363	0.683

The glances() function is an extension of glance() that can be used to look at many models simultaneously:

```
R> glances(spmod, lmod)
```

```
# A tibble: 2 x 10
```

	model	n	p	npar	value	AIC	AICc	logLik	deviance	pseudo.r.squared
	<chr>	<int>	<dbl>	<int>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	spmod	365	2	3	367.	373.	373.	-184.	363	0.683
2	lmod	365	2	1	634.	636.	636.	-317.	363.	0.671

Finally, the augment() function augments the original data with model diagnostics:

```

R> augment(spmod)

Simple feature collection with 365 features and 7 fields
Geometry type: POINT
Dimension:      XY
Bounding box:   xmin: -445884.1 ymin: 1929616 xmax: -383656.8 ymax: 2061414
Projected CRS:  NAD83 / Alaska Albers
# A tibble: 365 x 8
  log_Zn log_dist2road .fitted .resid   .hat .cooksd .std.resid
*   <dbl>         <dbl>   <dbl> <dbl>   <dbl>   <dbl>   <dbl>
1    7.33          2.68    8.26 -0.928 0.102   0.112   -1.48
2    7.38          2.68    8.26 -0.880 0.0101 0.000507 -0.316
3    7.58          2.54    8.34 -0.755 0.0170 0.000475 -0.236
4    7.63          2.97    8.09 -0.464 0.0137 0.000219  0.178
5    7.26          2.72    8.24 -0.977 0.0177 0.00515  -0.762
6    7.65          2.76    8.21 -0.568 0.0147 0.000929 -0.355
7    7.59          2.30    8.47 -0.886 0.0170 0.00802  -0.971
8    7.16          2.78    8.20 -1.05  0.0593 0.0492   -1.29
9    7.19          2.93    8.12 -0.926 0.00793 0.000451 -0.337
10   8.07          2.79    8.20 -0.123 0.0265 0.00396   0.547
# ... with 355 more rows, and 1 more variable: geometry <POINT [m]>

```

By default, only the columns of `data` used to fit the model are returned alongside the diagnostics. All original columns of `data` are returned by setting `drop` to `FALSE`. `augment()` is especially powerful when the data are an `sf` object because model diagnostics can be easily visualized spatially. For example, we could subset the augmented object so that it only includes observations whose standardized residuals have absolute values greater than some cutoff and then visualize them spatially. To learn more about the broom functions for spatial linear models, run `help("tidy.spm", "spmodel")`, `help("glance.spm", "spmodel")`, and `help("augment.spm", "spmodel")`.

An Areal Data Example

Next we use the `seal` data, an `sf` object that contains the log of the estimated harbor-seal trends from abundance data across polygons in Alaska, to provide an example of fitting a spatial linear model for areal data using `spautorm()`. We view the first few rows of `seal` by running

```

R> seal

Simple feature collection with 62 features and 1 field
Geometry type: POLYGON
Dimension:      XY
Bounding box:   xmin: 913618.8 ymin: 1007542 xmax: 1116002 ymax: 1145054
Projected CRS:  NAD83 / Alaska Albers
# A tibble: 62 x 2
  log_trend
  <dbl>
1    NA    ((1035002 1054710, 1035002 1054542, 1035002 1053542, 1035002 1052542, 1035002 1051542, 1035002 1050542, 1035002 1049542, 1035002 1048542, 1035002 1047542, 1035002 1046542, 1035002 1045542, 1035002 1044542, 1035002 1043542, 1035002 1042542, 1035002 1041542, 1035002 1040542, 1035002 1039542, 1035002 1038542, 1035002 1037542, 1035002 1036542, 1035002 1035542, 1035002 1034542, 1035002 1033542, 1035002 1032542, 1035002 1031542, 1035002 1030542, 1035002 1029542, 1035002 1028542, 1035002 1027542, 1035002 1026542, 1035002 1025542, 1035002 1024542, 1035002 1023542, 1035002 1022542, 1035002 1021542, 1035002 1020542, 1035002 1019542, 1035002 1018542, 1035002 1017542, 1035002 1016542, 1035002 1015542, 1035002 1014542, 1035002 1013542, 1035002 1012542, 1035002 1011542, 1035002 1010542, 1035002 1009542, 1035002 1008542, 1035002 1007542, 1035002 1006542, 1035002 1005542, 1035002 1004542, 1035002 1003542, 1035002 1002542, 1035002 1001542, 1035002 1000542, 1035002 999542, 1035002 998542, 1035002 997542, 1035002 996542, 1035002 995542, 1035002 994542, 1035002 993542, 1035002 992542, 1035002 991542, 1035002 990542, 1035002 989542, 1035002 988542, 1035002 987542, 1035002 986542, 1035002 985542, 1035002 984542, 1035002 983542, 1035002 982542, 1035002 981542, 1035002 980542, 1035002 979542, 1035002 978542, 1035002 977542, 1035002 976542, 1035002 975542, 1035002 974542, 1035002 973542, 1035002 972542, 1035002 971542, 1035002 970542, 1035002 969542, 1035002 968542, 1035002 967542, 1035002 966542, 1035002 965542, 1035002 964542, 1035002 963542, 1035002 962542, 1035002 961542, 1035002 960542, 1035002 959542, 1035002 958542, 1035002 957542, 1035002 956542, 1035002 955542, 1035002 954542, 1035002 953542, 1035002 952542, 1035002 951542, 1035002 950542, 1035002 949542, 1035002 948542, 1035002 947542, 1035002 946542, 1035002 945542, 1035002 944542, 1035002 943542, 1035002 942542, 1035002 941542, 1035002 940542, 1035002 939542, 1035002 938542, 1035002 937542, 1035002 936542, 1035002 935542, 1035002 934542, 1035002 933542, 1035002 932542, 1035002 931542, 1035002 930542, 1035002 929542, 1035002 928542, 1035002 927542, 1035002 926542, 1035002 925542, 1035002 924542, 1035002 923542, 1035002 922542, 1035002 921542, 1035002 920542, 1035002 919542, 1035002 918542, 1035002 917542, 1035002 916542, 1035002 915542, 1035002 914542, 1035002 913542, 1035002 912542, 1035002 911542, 1035002 910542, 1035002 909542, 1035002 908542, 1035002 907542, 1035002 906542, 1035002 905542, 1035002 904542, 1035002 903542, 1035002 902542, 1035002 901542, 1035002 900542, 1035002 899542, 1035002 898542, 1035002 897542, 1035002 896542, 1035002 895542, 1035002 894542, 1035002 893542, 1035002 892542, 1035002 891542, 1035002 890542, 1035002 889542, 1035002 888542, 1035002 887542, 1035002 886542, 1035002 885542, 1035002 884542, 1035002 883542, 1035002 882542, 1035002 881542, 1035002 880542, 1035002 879542, 1035002 878542, 1035002 877542, 1035002 876542, 1035002 875542, 1035002 874542, 1035002 873542, 1035002 872542, 1035002 871542, 1035002 870542, 1035002 869542, 1035002 868542, 1035002 867542, 1035002 866542, 1035002 865542, 1035002 864542, 1035002 863542, 1035002 862542, 1035002 861542, 1035002 860542, 1035002 859542, 1035002 858542, 1035002 857542, 1035002 856542, 1035002 855542, 1035002 854542, 1035002 853542, 1035002 852542, 1035002 851542, 1035002 850542, 1035002 849542, 1035002 848542, 1035002 847542, 1035002 846542, 1035002 845542, 1035002 844542, 1035002 843542, 1035002 842542, 1035002 841542, 1035002 840542, 1035002 839542, 1035002 838542, 1035002 837542, 1035002 836542, 1035002 835542, 1035002 834542, 1035002 833542, 1035002 832542, 1035002 831542, 1035002 830542, 1035002 829542, 1035002 828542, 1035002 827542, 1035002 826542, 1035002 825542, 1035002 824542, 1035002 823542, 1035002 822542, 1035002 821542, 1035002 820542, 1035002 819542, 1035002 818542, 1035002 817542, 1035002 816542, 1035002 815542, 1035002 814542, 1035002 813542, 1035002 812542, 1035002 811542, 1035002 810542, 1035002 809542, 1035002 808542, 1035002 807542, 1035002 806542, 1035002 805542, 1035002 804542, 1035002 803542, 1035002 802542, 1035002 801542, 1035002 800542, 1035002 799542, 1035002 798542, 1035002 797542, 1035002 796542, 1035002 795542, 1035002 794542, 1035002 793542, 1035002 792542, 1035002 791542, 1035002 790542, 1035002 789542, 1035002 788542, 1035002 787542, 1035002 786542, 1035002 785542, 1035002 784542, 1035002 783542, 1035002 782542, 1035002 781542, 1035002 780542, 1035002 779542, 1035002 778542, 1035002 777542, 1035002 776542, 1035002 775542, 1035002 774542, 1035002 773542, 1035002 772542, 1035002 771542, 1035002 770542, 1035002 769542, 1035002 768542, 1035002 767542, 1035002 766542, 1035002 765542, 1035002 764542, 1035002 763542, 1035002 762542, 1035002 761542, 1035002 760542, 1035002 759542, 1035002 758542, 1035002 757542, 1035002 756542, 1035002 755542, 1035002 754542, 1035002 753542, 1035002 752542, 1035002 751542, 1035002 750542, 1035002 749542, 1035002 748542, 1035002 747542, 1035002 746542, 1035002 745542, 1035002 744542, 1035002 743542, 1035002 742542, 1035002 741542, 1035002 740542, 1035002 739542, 1035002 738542, 1035002 737542, 1035002 736542, 1035002 735542, 1035002 734542, 1035002 733542, 1035002 732542, 1035002 731542, 1035002 730542, 1035002 729542, 1035002 728542, 1035002 727542, 1035002 726542, 1035002 725542, 1035002 724542, 1035002 723542, 1035002 722542, 1035002 721542, 1035002 720542, 1035002 719542, 1035002 718542, 1035002 717542, 1035002 716542, 1035002 715542, 1035002 714542, 1035002 713542, 1035002 712542, 1035002 711542, 1035002 710542, 1035002 709542, 1035002 708542, 1035002 707542, 1035002 706542, 1035002 705542, 1035002 704542, 1035002 703542, 1035002 702542, 1035002 701542, 1035002 700542, 1035002 699542, 1035002 698542, 1035002 697542, 1035002 696542, 1035002 695542, 1035002 694542, 1035002 693542, 1035002 692542, 1035002 691542, 1035002 690542, 1035002 689542, 1035002 688542, 1035002 687542, 1035002 686542, 1035002 685542, 1035002 684542, 1035002 683542, 1035002 682542, 1035002 681542, 1035002 680542, 1035002 679542, 1035002 678542, 1035002 677542, 1035002 676542, 1035002 675542, 1035002 674542, 1035002 673542, 1035002 672542, 1035002 671542, 1035002 670542, 1035002 669542, 1035002 668542, 1035002 667542, 1035002 666542, 1035002 665542, 1035002 664542, 1035002 663542, 1035002 662542, 1035002 661542, 1035002 660542, 1035002 659542, 1035002 658542, 1035002 657542, 1035002 656542, 1035002 655542, 1035002 654542, 1035002 653542, 1035002 652542, 1035002 651542, 1035002 650542, 1035002 649542, 1035002 648542, 1035002 647542, 1035002 646542, 1035002 645542, 1035002 644542, 1035002 643542, 1035002 642542, 1035002 641542, 1035002 640542, 1035002 639542, 1035002 638542, 1035002 637542, 1035002 636542, 1035002 635542, 1035002 634542, 1035002 633542, 1035002 632542, 1035002 631542, 1035002 630542, 1035002 629542, 1035002 628542, 1035002 627542, 1035002 626542, 1035002 625542, 1035002 624542, 1035002 623542, 1035002 622542, 1035002 621542, 1035002 620542, 1035002 619542, 1035002 618542, 1035002 617542, 1035002 616542, 1035002 615542, 1035002 614542, 1035002 613542, 1035002 612542, 1035002 611542, 1035002 610542, 1035002 609542, 1035002 608542, 1035002 607542, 1035002 606542, 1035002 605542, 1035002 604542, 1035002 603542, 1035002 602542, 1035002 601542, 1035002 600542, 1035002 599542, 1035002 598542, 1035002 597542, 1035002 596542, 1035002 595542, 1035002 594542, 1035002 593542, 1035002 592542, 1035002 591542, 1035002 590542, 1035002 589542, 1035002 588542, 1035002 587542, 1035002 586542, 1035002 585542, 1035002 584542, 1035002 583542, 1035002 582542, 1035002 581542, 1035002 580542, 1035002 579542, 1035002 578542, 1035002 577542, 1035002 576542, 1035002 575542, 1035002 574542, 1035002 573542, 1035002 572542, 1035002 571542, 1035002 570542, 1035002 569542, 1035002 568542, 1035002 567542, 1035002 566542, 1035002 565542, 1035002 564542, 1035002 563542, 1035002 562542, 1035002 561542, 1035002 560542, 1035002 559542, 1035002 558542, 1035002 557542, 1035002 556542, 1035002 555542, 1035002 554542, 1035002 553542, 1035002 552542, 1035002 551542, 1035002 550542, 1035002 549542, 1035002 548542, 1035002 547542, 1035002 546542, 1035002 545542, 1035002 544542, 1035002 543542, 1035002 542542, 1035002 541542, 1035002 540542, 1035002 539542, 1035002 538542, 1035002 537542, 1035002 536542, 1035002 535542, 1035002 534542, 1035002 533542, 1035002 532542, 1035002 531542, 1035002 530542, 1035002 529542, 1035002 528542, 1035002 527542, 1035002 526542, 1035002 525542, 1035002 524542, 1035002 523542, 1035002 522542, 1035002 521542, 1035002 520542, 1035002 519542, 1035002 518542, 1035002 517542, 1035002 516542, 1035002 515542, 1035002 514542, 1035002 513542, 1035002 512542, 1035002 511542, 1035002 510542, 1035002 509542, 1035002 508542, 1035002 507542, 1035002 506542, 1035002 505542, 1035002 504542, 1035002 503542, 1035002 502542, 1035002 501542, 1035002 500542, 1035002 499542, 1035002 498542, 1035002 497542, 1035002 496542, 1035002 495542, 1035002 494542, 1035002 493542, 1035002 492542, 1035002 491542, 1035002 490542, 1035002 489542, 1035002 488542, 1035002 487542, 1035002 486542, 1035002 485542, 1035002 484542, 1035002 483542, 1035002 482542, 1035002 481542, 1035002 480542, 1035002 479542, 1035002 478542, 1035002 477542, 1035002 476542, 1035002 475542, 1035002 474542, 1035002 473542, 1035002 472542, 1035002 471542, 1035002 470542, 1035002 469542, 1035002 468542, 1035002 467542, 1035002 466542, 1035002 465542, 1035002 464542, 1035002 463542, 1035002 462542, 1035002 461542, 1035002 460542, 1035002 459542, 1035002 458542, 1035002 457542, 1035002 456542, 1035002 455542, 1035002 454542, 1035002 453542, 1035002 452542, 1035002 451542, 1035002 450542, 1035002 449542, 1035002 448542, 1035002 447542, 1035002 446542, 1035002 445542, 1035002 444542, 1035002 443542, 1035002 442542, 1035002 441542, 1035002 440542, 1035002 439542, 1035002 438542, 1035002 437542, 1035002 436542, 1035002 435542, 1035002 434542, 1035002 433542, 1035002 432542, 1035002 431542, 1035002 430542, 1035002 429542, 1035002 428542, 1035002 427542, 1035002 426542, 1035002 425542, 1035002 424542, 1035002 423542, 1035002 422542, 1035002 421542, 1035002 420542, 1035002 419542, 1035002 418542, 1035002 417542, 1035002 416542, 1035002 415542, 1035002 414542, 1035002 413542, 1035002 412542, 1035002 411542, 1035002 410542, 1035002 409542, 1035002 408542, 1035002 407542, 1035002 406542, 1035002 405542, 1035002 404542, 1035002 403542, 1035002 402542, 1035002 401542, 1035002 400542, 1035002 399542, 1035002 398542, 1035002 397542, 1035002 396542, 1035002 395542, 1035002 394542, 1035002 393542, 1035002 392542, 1035002 391542, 1035002 390542, 1035002 389542, 1035002 388542, 1035002 387542, 1035002 386542, 1035002 385542, 1035002 384542, 1035002 383542, 1035002 382542, 1035002 381542, 1035002 380542, 1035002 379542, 1035002 378542, 1035002 377542, 1035002 376542, 1035002 375542, 1035002 374542, 1035002 373542, 1035002 372542, 1035002 371542, 1035002 370542, 1035002 369542, 1035002 368542, 1035002 367542, 1035002 366542, 1035002 365542, 1035002 364542
```

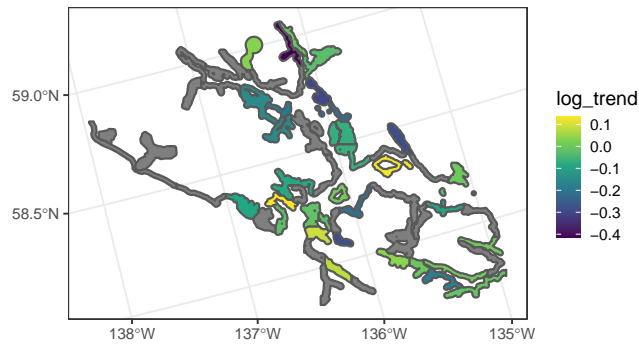


Fig 3. Distribution of log seal trends in the seal data. Polygons are gray if seal trends are missing.

```

6    0.0872 ((1026035 1044623, 1026037 1044605, 1026072 1044610, 1026083 10446
7   -0.266  ((1100345 1060709, 1100287 1060706, 1100228 1060706, 1100170 10607
8    0.0743 ((1030247 1029637, 1030248 1029637, 1030265 1029642, 1030328 10296
9   NA      ((1043093 1020553, 1043097 1020550, 1043101 1020550, 1043166 10205
10  -0.00961 ((1116002 1024542, 1116002 1023542, 1116002 1022542, 1116002 10215
# ... with 52 more rows

```

We can learn more about the data by running `help("seal", "spmodel")`.

We can visualize the distribution of log seal trends in the `seal` data (Figure 3) by running

```
R> ggplot(seal, aes(fill = log_trend)) +
+   geom_sf(size = 0.75) +
+   scale_fill_viridis_c() +
+   theme_bw(base_size = 14)
```

The gray polygons denote areas where the log trend is missing. These missing areas need to be kept in the data while fitting the model to preserve the overall neighborhood structure.

We estimate parameters of a spatial autoregressive model for log seal trends (`log_trend`) using an intercept-only model with a conditional autoregressive (CAR) spatial covariance by running

```
R> sealmod <- spautor(log_trend ~ 1, seal, spcov_type = "car")
```

If a weight matrix is not provided to `spautor()`, it is calculated internally using queen contiguity. Recall that queen contiguity defines two observations as neighbors if they share at least one common boundary. If at least one observation has no neighbors, the `extra` parameter is estimated, which quantifies variability among observations without neighbors. By default, `spautor()` uses row standardization [19] and assumes an independent error variance (`ie`) of zero.

We tidy and glance at the fitted model and augment the data by running

```
R> tidy(sealmod)
```

```

# A tibble: 1 x 5
  term          estimate std.error statistic p.value
<chr>         <dbl>     <dbl>     <dbl>   <dbl>
1 (Intercept) -0.0710     0.0250     -2.85 0.00443

R> glance(sealmod)

# A tibble: 1 x 9
  n     p npar value   AIC  AICc logLik deviance pseudo.r.squared
<int> <dbl> <int> <dbl> <dbl> <dbl> <dbl>   <dbl>         <dbl>
1    34     1     3 -36.9 -30.9 -30.1  18.4    32.9             0

R> augment(sealmod)

Simple feature collection with 34 features and 6 fields
Geometry type: POLYGON
Dimension: XY
Bounding box: xmin: 980001.5 ymin: 1010815 xmax: 1116002 ymax: 1145054
Projected CRS: NAD83 / Alaska Albers
# A tibble: 34 x 7
  log_trend .fitted .resid .hat .cooksd .std.resid
*   <dbl>   <dbl>   <dbl> <dbl> <dbl>   <dbl>
1 -0.282 -0.0710 -0.211 0.0179 0.0233 -1.14 ((1037002 1039492, 10370
2 -0.00121 -0.0710 0.0698 0.0699 0.0412 0.767 ((1070158 1030216, 10701
3 0.0354 -0.0710 0.106 0.0218 0.0109 0.705 ((1054906 1034826, 10549
4 -0.0160 -0.0710 0.0550 0.0343 0.00633 0.430 ((1025142 1056940, 10251
5 0.0872 -0.0710 0.158 0.0229 0.0299 1.14 ((1026035 1044623, 10260
6 -0.266 -0.0710 -0.195 0.0280 0.0493 -1.33 ((1100345 1060709, 11002
7 0.0743 -0.0710 0.145 0.0480 0.0818 1.30 ((1030247 1029637, 10302
8 -0.00961 -0.0710 0.0614 0.0143 0.00123 0.293 ((1116002 1024542, 11160
9 -0.182 -0.0710 -0.111 0.0131 0.0155 -1.09 ((1079864 1025088, 10798
10 0.00351 -0.0710 0.0745 0.0340 0.0107 0.561 ((1110363 1037056, 11103
# ... with 24 more rows

```

Prediction

In this section, we show how to use `predict()` to perform spatial prediction (also called Kriging) in `splmodel`. We will fit a model using the point-referenced `sulfate` data, an `sf` object that contains sulfate measurements in the conterminous United States, and make predictions for each location in the point-referenced `sulfate_preds` data, an `sf` object that contains locations in the conterminous United States at which to predict sulfate.

We first visualize the distribution of the sulfate data (Figure 4, left) by running

```

R> ggplot(sulfate, aes(color = sulfate)) +
+   geom_sf(size = 2.5) +
+   scale_color_viridis_c(limits = c(0, 45)) +
+   theme_gray(base_size = 18)

```

We then fit a spatial linear model for sulfate using an intercept-only model with a spherical spatial covariance function by running

```

R> sulfmod <- splm(sulfate ~ 1, sulfate, spcov_type = "spherical")

```

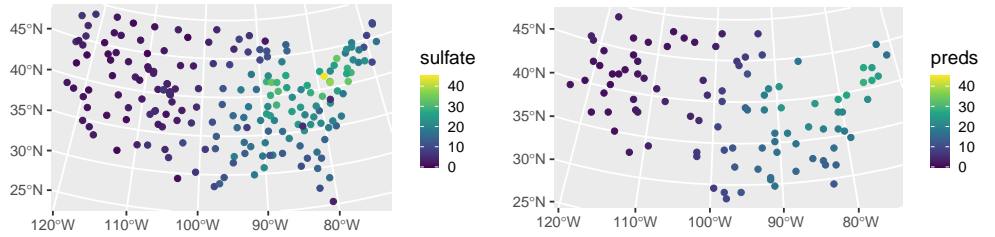


Fig 4. Distribution of observed sulfate (left) and sulfate predictions (right) in the conterminous United States.

Then we obtain best linear unbiased predictions (Kriging predictions) using `predict()`, where the `newdata` argument contains the locations at which to predict, storing them as a new variable in `sulfate_preds` called `preds`:

```
R> sulfate_preds$preds <- predict(sulfmod, newdata = sulfate_preds)
```

We can then visualize the model predictions (Figure 4, right) by running

```
R> sulfate_preds$preds <- predict(sulfmod, newdata = sulfate_preds)
R> ggplot(sulfate_preds, aes(color = preds)) +
+   geom_sf(size = 2.5) +
+   scale_color_viridis_c(limits = c(0, 45)) +
+   theme_gray(base_size = 18)
```

Before making predictions, it is important to properly specify the `newdata` object. If explanatory variables were used to fit the model, the same explanatory variables must be included in `newdata` with the same names they have in `data`. If `data` is a `data.frame`, coordinates must be included in `newdata` with the same names as they have in `data`. If `data` is an `sf` object, coordinates must be included in `newdata` with the same geometry name as they have in `data`. When using projected coordinates, the projection for `newdata` should be the same as the projection for `data`.

Prediction standard errors are returned by setting the `se.fit` argument to `TRUE`:

```
R> predict(sulfmod, newdata = sulfate_preds, se.fit = TRUE)
```

The `interval` argument determines the type of interval returned. If `interval` is `"none"` (the default), no intervals are returned. If `interval` is `"prediction"`, `100 * level%` prediction intervals are returned (the default is 95% prediction intervals):

```
R> predict(sulfmod, newdata = sulfate_preds, interval = "prediction")
```

If `interval` is `"confidence"`, the predictions are instead the estimated mean given each observation's explanatory variable values. The corresponding `100 * level%` confidence intervals are returned:

```
R> predict(sulfmod, newdata = sulfate_preds, interval = "confidence")
```

Previously we used the `augment()` function to augment `data` with model diagnostics. We can also use `augment()` to augment `newdata` with predictions, standard errors, and intervals. We remove the model predictions from `sulfate_preds` before showing how `augment()` is used to obtain the same predictions by running

```

R> sulfate_preds$preds <- NULL
470

We then view the first few rows of sulfate_preds augmented with a 90% prediction
471 interval by running
472

R> augment(sulfmod, newdata = sulfate_preds, interval = "prediction", level = 0.9)
473

Simple feature collection with 100 features and 3 fields
474
Geometry type: POINT
475
Dimension: XY
476
Bounding box: xmin: -2283774 ymin: 582930.5 xmax: 1985906 ymax: 3037173
477
Projected CRS: NAD83 / Conus Albers
478
# A tibble: 100 x 4
479
  .fitted .lower .upper geometry
  <dbl> <dbl> <dbl> <POINT [m]>
480
1 1.40 -5.33 8.14 (-1771413 1752976)
481
2 24.5 18.2 30.8 (1018112 1867127)
482
3 8.99 2.36 15.6 (-291256.8 1553212)
483
4 16.4 9.92 23.0 (1274293 1267835)
484
5 4.91 -1.56 11.4 (-547437.6 1638825)
485
6 26.7 20.4 33.0 (1445080 1981278)
486
7 3.00 -3.65 9.66 (-1629090 3037173)
487
8 14.3 7.97 20.6 (1302757 1039534)
488
9 1.49 -5.08 8.06 (-1429838 2523494)
489
10 14.4 7.97 20.8 (1131970 1096609)
490
# ... with 90 more rows
492

Here .fitted represents the predictions.
493

An alternative (but equivalent) approach can be used for model fitting and
494 prediction that circumvents the need to keep data and newdata as separate objects.
495 Suppose that observations requiring prediction are stored in data as missing (NA) values.
496 We can add a column of missing values to sulfate_preds and then bind it together
497 with sulfate by running
498

R> sulfate_preds$sulfate <- NA
499
R> sulfate_with_NA <- rbind(sulfate, sulfate_preds)
500

We can then fit a spatial linear model by running
501

R> sulfmod_with_NA <- splm(sulfate ~ 1, sulfate_with_NA, "spherical")
502

The missing values are ignored for model-fitting but stored in sulfmod_with_NA as
503 newdata:
504

R> sulfmod_with_NA$newdata
505

Simple feature collection with 100 features and 1 field
506
Geometry type: POINT
507
Dimension: XY
508
Bounding box: xmin: -2283774 ymin: 582930.5 xmax: 1985906 ymax: 3037173
509
Projected CRS: NAD83 / Conus Albers
510
First 10 features:
511
  sulfate geometry
512
198 NA POINT (-1771413 1752976)
513
199 NA POINT (1018112 1867127)
514

```

```

200      NA POINT (-291256.8 1553212)
201      NA  POINT (1274293 1267835)
202      NA POINT (-547437.6 1638825)
203      NA  POINT (1445080 1981278)
204      NA  POINT (-1629090 3037173)
205      NA  POINT (1302757 1039534)
206      NA  POINT (-1429838 2523494)
207      NA  POINT (1131970 1096609)

```

We can then predict the missing values by running

```
R> predict(sulfmod_with_NA)
```

The call to `predict()` finds in `sulfmod_with_NA` the `newdata` object and is equivalent to

```
R> predict(sulfmod_with_NA, newdata = sulfmod_with_NA$newdata)
```

We can also use `augment()` to make the predictions on the data set with missing values by running

```
R> augment(sulfmod_with_NA, newdata = sulfmod_with_NA$newdata)
```

```

Simple feature collection with 100 features and 2 fields
Geometry type: POINT
Dimension:     XY
Bounding box:  xmin: -2283774 ymin: 582930.5 xmax: 1985906 ymax: 3037173
Projected CRS: NAD83 / Conus Albers
# A tibble: 100 x 3
  sulfate .fitted geometry
*   <dbl>   <dbl>   <POINT [m]>
1     NA    1.40 (-1771413 1752976)
2     NA   24.5  (1018112 1867127)
3     NA    8.99 (-291256.8 1553212)
4     NA   16.4  (1274293 1267835)
5     NA    4.91 (-547437.6 1638825)
6     NA   26.7  (1445080 1981278)
7     NA    3.00 (-1629090 3037173)
8     NA   14.3  (1302757 1039534)
9     NA    1.49 (-1429838 2523494)
10    NA   14.4  (1131970 1096609)
# ... with 90 more rows

```

Unlike `predict()`, `augment()` explicitly requires the `newdata` argument be specified in order to obtain predictions. Omitting `newdata` (e.g., running `augment(sulfmod_with_NA)`) returns model diagnostics, not predictions.

For areal data models fit with `spautor()`, predictions cannot be computed at locations that were not incorporated in the neighborhood structure used to fit the model. Thus, predictions are only possible for observations in `data` whose response values are missing (NA), as their locations are incorporated into the neighborhood structure. For example, we make predictions of log seal trends at the missing polygons from Figure 3 by running

```
R> predict(sealmod)
```



```

$ is_known
  ie
TRUE

attr(,"class")
[1] "exponential"

```

The `init` output shows that the `ie` parameter has an initial value of zero that is assumed to be known. Next the model is fit:

```
R> spmod_red <- splm(log_Zn ~ log_dist2road, moss, spcov_initial = init)
```

Notice that because the `spcov_initial` object contains information about the spatial covariance type, the `spcov_type` argument is not required when `spcov_initial` is provided. We can use `glances()` to glance at both models:

```
R> glances(spmod, spmod_red)
```

```
# A tibble: 2 x 10
```

model	n	p	npar	value	AIC	AICc	logLik	deviance	pseudo.r.squared
<chr>	<int>	<dbl>	<int>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1 spmod	365	2	3	367.	373.	373.	-184.	363	0.68
2 spmod_red	365	2	2	378.	382.	382.	-189.	374.	0.70

The lower AIC and AICc of the full model compared to the reduced model indicates that the independent random error variance is important to the model. A likelihood ratio test comparing the full and reduced models is also possible using `anova()`.

Another application of fixing spatial covariance parameters involves calculating their profile likelihood confidence intervals [33]. Before calculating a profile likelihood confidence interval for Θ_i , the i th element of a general parameter vector Θ , it is necessary to obtain $-2\ell(\hat{\Theta})$, minus twice the log-likelihood evaluated at the estimated parameter vector, $\hat{\Theta}$. Then a $(1 - \alpha)\%$ profile likelihood confidence interval is then the set of values for Θ_i such that $2\ell(\hat{\Theta}) - 2\ell(\hat{\Theta}_{-i}) \leq \chi^2_{1,1-\alpha}$, where $\ell(\hat{\Theta}_{-i})$ is the value of the log-likelihood maximized after fixing Θ_i and optimizing over the remaining parameters, Θ_{-i} , and $\chi^2_{1,1-\alpha}$ is the $1 - \alpha$ quantile of a chi-squared distribution with one degree of freedom. The result follows from inverting a likelihood ratio test comparing the full model to a reduced model that fixes the value of Θ_i . Because computing profile likelihood confidence intervals requires refitting the model many times for different fixed values of Θ_i , it can be computationally intensive. This approach can be generalized to yield joint profile likelihood confidence intervals cases when i has dimension greater than one.

Random Effects

Non-spatial random effects incorporate additional sources of variability into model fitting. They are accommodated in `spmodel` using similar syntax as for random effects in the `nlme` [27] and `lme4` [34] **R** packages. Random effects are specified via a formula passed to the `random` argument. Next, we show two examples that incorporate random effects into the spatial linear model using the `moss` data.

The first example explores random intercepts for the `sample` variable. The `sample` variable indexes each unique location, which can have replicate observations due to field duplicates (`field_dup`) and lab replicates (`lab_rep`). There are 365 observations in `moss` at 318 unique locations, which means that 47 observations in `moss` are either field

duplicates or lab replicates. It is likely that the repeated observations at a location are correlated with one another. We can incorporate this repeated-observation correlation by creating a random intercept for each level of `sample`. We model the random intercepts for `sample` by running

```
R> rand1 <- splm(
+   log_Zn ~ log_dist2road,
+   moss,
+   spcov_type = "exponential",
+   random = ~ sample
+ )
```

Note that `~ sample` is shorthand for `(1 | sample)`, which is more explicit notation that indicates random intercepts for each level of `sample`.

The second example adds a random intercept for `year`, which creates extra correlation for observations within a year. It also adds a random slope for `log_dist2road` within `year`, which lets the effect of `log_dist2road` vary between years. We fit this model by running

```
R> rand2 <- splm(
+   log_Zn ~ log_dist2road,
+   moss,
+   spcov_type = "exponential",
+   random = ~ sample + (log_dist2road | year)
+ )
```

Note that `sample + (log_dist2road | year)` is shorthand for `(1 | sample) + (log_dist2road | year)`. If only random slopes within a year are desired (and no random intercepts), a `- 1` is given to the relevant portion of the formula: `(log_dist2road - 1 | year)`. When there is more than one term in `random`, each term must be surrounded by parentheses (recall that the random intercept shorthand automatically includes relevant parentheses).

We glance at all three models by running

```
R> glances(spmo, rand1, rand2)
```

```
# A tibble: 3 x 10
```

	model	n	p	npar	value	AIC	AICc	logLik	deviance	pseudo.r.squared
	<chr>	<int>	<dbl>	<int>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	rand2	365	2	6	190.	202.	202.	-94.9	363.	0.215
2	rand1	365	2	4	335.	343.	343.	-168.	363	0.661
3	spmo	365	2	3	367.	373.	373.	-184.	363	0.683

`rand2` has the lowest AIC and AICc.

It is possible to fix random effect variances using the `randcov_initial` argument, and `randcov_initial` can also be used to set initial values for optimization.

Partition Factors

A partition factor is a variable that allows observations to be uncorrelated when they are from different levels of the partition factor. Partition factors are specified in `spmodel` by providing a formula with a single variable to the `partition_factor` argument. Suppose that for the `moss` data, we would like observations in different years (`year`) to be uncorrelated. We fit a model that treats `year` as a partition factor by running

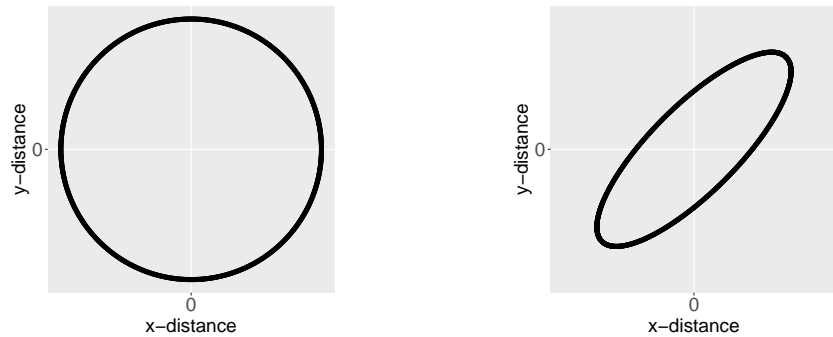


Fig 5. Ellipses for an isotropic (left) and anisotropic (right) covariance function centered at the origin. The black outline of each ellipse is a level curve of equal correlation.

```
R> part <- splm(
+   log_Zn ~ log_dist2road,
+   moss,
+   spcov_type = "exponential",
+   partition_factor = ~ year
+ )
```

Anisotropy

A spatial covariance function for point-referenced data is isotropic if it behaves similarly in all directions (i.e., is independent of direction) as a function of distance. An anisotropic covariance function does not behave similarly in all directions as a function of distance. Consider the spatial covariance imposed by an eastward-moving wind pattern. A one-unit distance in the x-direction likely means something different than a one-unit distance in the y-direction. Figure 5 shows ellipses for an isotropic and anisotropic covariance function centered at the origin (a distance of zero).

The black outline of each ellipse is a level curve of equal correlation. The left ellipse (a circle) represents an isotropic covariance function. The distance at which the correlation between two observations lays on the level curve is the same in all directions. The right ellipse represents an anisotropic covariance function. The distance at which the correlation between two observations lays on the level curve is different in different directions.

Accounting for anisotropy involves a rotation and scaling of the x-coordinates and y-coordinates such that the spatial covariance function that uses these transformed distances is isotropic. We use the `anisotropy` argument to `splm()` to fit a model with anisotropy by running

```
R> spmod_anis <- splm(
+   log_Zn ~ log_dist2road,
+   moss,
+   spcov_type = "exponential",
+   anisotropy = TRUE
+ )
R> summary(spmod_anis)
```

Call:

```
splm(formula = log_Zn ~ log_dist2road, data = moss, spcov_type = "exponential",
      anisotropy = TRUE)
```

```

Residuals:
    Min       1Q   Median       3Q      Max
-2.5279 -1.2239 -0.7202 -0.1921  1.1659

Coefficients (fixed):
              Estimate Std. Error z value Pr(>|z|)
(Intercept)   9.54798    0.22291   42.83  <2e-16 ***
log_dist2road -0.54601    0.01855  -29.44  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Pseudo R-squared: 0.7048

Coefficients ( spatial covariance):
           de           ie        range        rotate        scale
3.561e-01 6.812e-02 8.732e+03 2.435e+00 4.753e-01
attr(,"class")
[1] "exponential"

```

The **rotate** parameter is between zero and π radians and represents the angle of a clockwise rotation of the ellipse such that the major axis of the ellipse is the new x-axis and the minor axis of the ellipse is the new y-axis. The **scale** parameter is between zero and one and represents the ratio of the distance between the origin and the edge of the ellipse along the minor axis to the distance between the origin and the edge of the ellipse along the major axis. Figure 6 shows a transformation that turns an anisotropic ellipse into an isotropic one (i.e., a circle). This transformation requires rotating the coordinates clockwise by **rotate** and then scaling them the reciprocal of **scale**. The transformed coordinates are then used instead of the original coordinates to compute distances and spatial covariances.

Note that specifying an initial value for **rotate** that is different from zero, specifying an initial value for **scale** that is different from one, or assuming either **rotate** or **scale** are unknown in **spcov_initial** will cause **splm()** to fit a model with anisotropy (and will override **anisotropy = FALSE**). Estimating anisotropy parameters is only possible for maximum likelihood and restricted maximum likelihood estimation, but fixed anisotropy parameters can be accommodated for semivariogram weighted least squares or semivariogram composite likelihood estimation. Also note that anisotropy is not relevant for areal data because the spatial covariance function depends on a neighborhood structure instead of distances between points.

Simulating Spatial Data

The **sprnorm()** function is used to simulate normal (Gaussian) spatial data. To use **sprnorm()**, the **spcov_params()** function is used to create an **spcov_params** object. The **spcov_params()** function requires the spatial covariance type and parameter values. We create an **spcov_params** object by running

```
R> sim_params <- spcov_params("exponential", de = 5, ie = 1, range = 0.5)
```

We set a reproducible seed and then simulate data at 3000 random locations in the unit square using the spatial covariance parameters in **sim_params** by running

```
R> set.seed(0)
R> n <- 3000
```



Fig 6. A visual representation of the anisotropy transformation. In the left figure, the first step is to rotate the anisotropic ellipse clockwise by the `rotate` parameter (here `rotate` is 0.75 radians or 135 degrees). In the middle figure, the second step is to scale the y axis by the reciprocal of the `scale` parameter (here `scale` is 0.5). In the right figure, the anisotropic ellipse has been transformed into an isotropic one (i.e., a circle). The transformed coordinates are then used instead of the original coordinates to compute distances and spatial covariances.

```
R> x <- runif(n) 778
R> y <- runif(n) 779
R> sim_coords <- tibble::tibble(x, y) 780
R> sim_response <- sprnorm(sim_params, data = sim_coords, xcoord = x, ycoord = y) 781
R> sim_data <- tibble::tibble(sim_coords, sim_response) 782
```

We can visualize the simulated data (Figure 7, left) by running

```
R> ggplot(sim_data, aes(x = x, y = y, color = sim_response)) + 784
+   geom_point(size = 1.5) + 785
+   scale_color_viridis_c(limits = c(-7, 7)) + 786
+   theme_gray(base_size = 18) 787
```

There is noticeable spatial patterning in the response variable (`sim_response`). The default mean in `sprnorm()` is zero for all observations, though a mean vector can be provided using the `mean` argument. The default number of samples generated in `sprnorm()` is one, though this can be changed using the `samples` argument. Because `sim_data` is a `tibble` (`data.frame`) and not an `sf` object, the columns in `sim_data` representing the x-coordinates and y-coordinates must be provided to `sprnorm()`.

Note that the output from `coef(object, type = "spcov")` is a `spcov_params` object. This is useful we want to simulate data given the estimated spatial covariance parameters from a fitted model. Random effects are incorporated into simulation via the `randcov_params` argument.

Big Data

The computational cost associated with model fitting is exponential in the sample size for all estimation methods. For maximum likelihood and restricted maximum likelihood, the computational cost of estimating θ is cubic. For semivariogram weighted least squares and semivariogram composite likelihood, the computational cost of estimating θ

is quadratic. The computational cost associated with estimating β and prediction is cubic in the model-fitting sample size, regardless of estimation method. Typically, samples sizes approaching 10,000 make the computational cost of model fitting and prediction infeasible, which necessitates the use of big data methods. `spmodel` offers big data methods for model fitting of point-referenced data via the `local` argument to `splm()`. The method is capable of quickly fitting models with hundreds of thousands to millions of observations. Because of the neighborhood structure of areal data, the big data methods used for point-referenced data do not apply to areal data. Thus, there is no big data method for areal data or `local` argument to `spautorm()`, so model fitting sample sizes cannot be too large

`spmodel` offers big data methods for prediction of point-referenced data or areal data via the `local` argument to `predict()`, capable of quickly predicting hundreds of thousands to millions of observations rather quickly. To show how to use `spmodel` for big data estimation and prediction, we use the `sim_data` data from Section . Because `sim_data` is a `tibble` (`data.frame`) and not an `sf` object, the columns in `data` representing the x-coordinates and y-coordinates must be explicitly provided to `splm()`.

Model-fitting

`spmodel` uses a “local indexing” approximation for big data model fitting of point-referenced data. Observations are first assigned an index. Then for the purposes of model fitting, observations with different indexes are assumed uncorrelated. Assuming observations with different indexes are uncorrelated induces sparsity in the covariance matrix, which greatly reduces the computational time of operations that involve the covariance matrix.

The `local` argument to `splm()` controls the big data options. `local` is a list with several arguments. The arguments to the `local` list control the method used to assign the indexes, the number of observations with the same index, the number of unique indexes, variance adjustments to the covariance matrix of $\hat{\beta}$, whether or not to use parallel processing, and if parallel processing is used, the number of cores.

The simplest way to accommodate big data is to set `local` to `TRUE`. This is shorthand for `local = list(method = "random", size = 50, var_adjust = "theoretical", parallel = FALSE)`, which randomly assigns observations to groups, ensures each group has approximately 50 observations, uses the theoretically-correct variance adjustment, and does not use parallel processing.

```
R> local1 <- splm(sim_response ~ 1, sim_data, spcov_type = "exponential",
+               xcoord = x, ycoord = y, local = TRUE)
R> summary(local1)
```

```
Call:
splm(formula = sim_response ~ 1, data = sim_data, spcov_type = "exponential",
      xcoord = x, ycoord = y, local = TRUE)
```

Residuals:

Min	1Q	Median	3Q	Max
-5.0356	-1.3514	-0.1468	1.2842	6.5381

Coefficients (fixed):

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-1.021	0.699	-1.46	0.144

Coefficients (spatial covariance):

```

      de      ie range
2.8724 0.9735 0.2644

```

Instead of using `local = TRUE`, we can explicitly set `local`. For example, we can fit a model using k-means clustering [35] on the x-coordinates and y-coordinates to create 60 groups (clusters), use the pooled variance adjustment, and use parallel processing with two cores by running

```

R> local2_list <- list(method = "kmeans", groups = 60, var_adjust = "pooled",
+                      parallel = TRUE, ncores = 2)
R> local2 <- splm(sim_response ~ 1, sim_data, spcov_type = "exponential",
+               xcoord = x, ycoord = y, local = local2_list)
R> summary(local2)

```

```

Call:
splm(formula = sim_response ~ 1, data = sim_data, spcov_type = "exponential",
      xcoord = x, ycoord = y, local = local2_list)

```

```

Residuals:
      Min       1Q   Median       3Q      Max
-4.98801 -1.30386 -0.09927  1.33176  6.58567

```

```

Coefficients (fixed):
              Estimate Std. Error z value Pr(>|z|)
(Intercept)  -1.0683      0.1759  -6.073 1.25e-09 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

Coefficients ( spatial covariance):
      de      ie range
2.5434 0.9907 0.2312

```

Likelihood-based statistics like `AIC()`, `AICc()`, `logLik()`, and `deviance()` should not be used to compare a model fit with a big data approximation to a model fit without a big data approximation, as the two approaches maximize different likelihoods.

Prediction

For point-referenced data, `splm` uses a “local neighborhood” approximation for big data prediction. Each prediction is computed using a subset of the observed data instead of all of the observed data. Before further discussing big data prediction, we simulate 1000 locations in the unit square requiring prediction:

```

R> n_pred <- 1000
R> x <- runif(n_pred)
R> y <- runif(n_pred)
R> sim_preds <- tibble::tibble(x = x, y = y)

```

The `local` argument to `predict()` controls the big data options. `local` is a list with several arguments. The arguments to the `local` list control the method used to subset the observed data, the number of observations in each subset, whether or not to use parallel processing, and if parallel processing is used, the number of cores.

The simplest way to accommodate big data prediction is to set `local` to `TRUE`. This is shorthand for `local = list(method = "covariance", size = 50, parallel =`

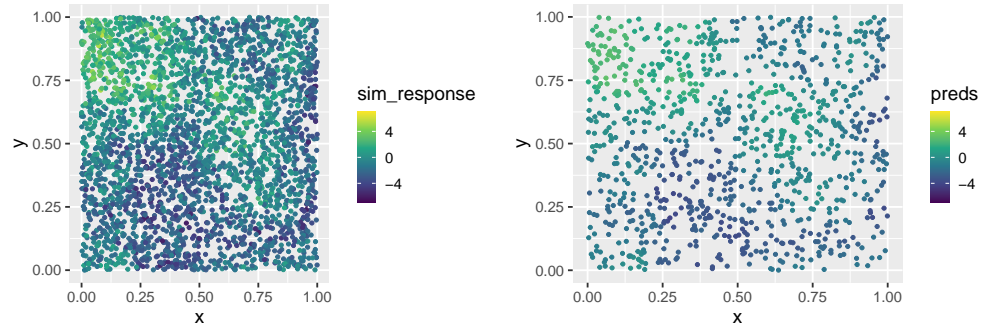


Fig 7. Observed data and big data predictions at unobserved locations. In the left figure, spatial data are simulated in the unit square. A spatial linear model is fit using the default big data approximation for model-fitting. In the right figure, predictions are made using the fitted model and the default big data approximation for prediction.

`FALSE`), which implies that, for each location requiring prediction, only the 50 observations in the data most correlated with it are used in the computation and parallel processing is not used. Using the `local1` fitted model, we store these predictions as a variable called `preds` in the `sim_preds` data by running

```
R> sim_preds$preds <- predict(local1, newdata = sim_preds, local = TRUE)
```

The predictions are visualized (Figure 7, right) by running

```
R> ggplot(sim_preds, aes(x = x, y = y, color = preds)) +
+   geom_point(size = 1.5) +
+   scale_color_viridis_c(limits = c(-7, 7)) +
+   theme_gray(base_size = 18)
```

They display a similar pattern as the observed data.

Instead of using `local = TRUE`, we can explicitly set `local`:

```
R> pred_list <- list(method = "distance", size = 30, parallel = TRUE, ncores = 1)
R> predict(local1, newdata = sim_preds, local = pred_list)
```

This code implies that uniquely for each location requiring prediction, only the 30 observations in the data closest to it (in terms of Euclidean distance) are used in the computation and parallel processing is used with two cores.

For areal data, no local neighborhood approximation exists because of the data's underlying neighborhood structure. Thus, all of the data must be used to compute predictions and by consequence, `method` and `distance` are not components of the `local` list. The only components of the `local` list for areal data are `parallel` and `ncores`.

Discussion

`splm` is a novel, relevant contribution used to fit, summarize, and predict for a variety of spatial statistical models. Spatial linear models for point-referenced data (i.e., geostatistical models) are fit using the `splm()` function. Spatial linear models for areal data (i.e., spatial autoregressive models) are fit using the `spautorm()` function. Both functions share use a common framework and syntax structure. Several model-fit statistics and diagnostics are available. The broom functions `tidy()` and `glance()` are used to tidy and glance at a fitted model. The broom function `augment()` is used to

augment **data** with model diagnostics and augment **newdata** with predictions. Several
advanced features are available to accommodate fixed covariance parameter values,
random effects, partition factors, anisotropy, simulated data, and big data
approximations for model fitting and prediction.

We appreciate feedback from users regarding **spmodel**. To learn more about how to
provide feedback or contribute to **spmodel**, please visit our GitHub repository at
<https://github.com/USEPA/spmodel>.

Acknowledgements

The views expressed in this manuscript are those of the authors and do not necessarily
represent the views or policies of the U.S. Environmental Protection Agency or the
National Oceanic and Atmospheric Administration. Any mention of trade names,
products, or services does not imply an endorsement by the U.S. government, the U.S.
Environmental Protection Agency, or the National Oceanic and Atmospheric
Administration. The U.S. Environmental Protection Agency and the National Oceanic
and Atmospheric Administration do not endorse any commercial products, services or
enterprises.

References

1. Dumelle M, Kincaid TM, Olsen AR, Weber MH. Spsurvey: Spatial sampling
design and analysis. 2022. Available: <https://CRAN.R-project.org/package=spsurvey>
2. Dumelle M, Higham M, Ver Hoef JM, Olsen AR, Madsen L. A comparison of
design-based and model-based approaches for finite population spatial sampling
and inference. *Methods in Ecology and Evolution*. 2022.
3. Cressie N. *Statistics for spatial data*. John Wiley & Sons; 1993.
4. Nychka D, Furrer R, Paige J, Sain S. *Fields: Tools for spatial data*. Boulder,
CO, USA: University Corporation for Atmospheric Research; 2021. Available:
<https://github.com/dnychka/fieldsRPackage>
5. Zammit-Mangion A, Cressie N. FRK: An r package for spatial and spatio-temporal
prediction with large datasets. *Journal of Statistical Software*. 2021;98: 1–48.
6. Ribeiro Jr PJ, Diggle P, Christensen O, Schlather M, Bivand R, Ripley B. *geoR:*
Analysis of geostatistical data. 2022. Available: <https://CRAN.R-project.org/package=geoR>
7. Guinness J, Katzfuss M, Fahmy Y. GpGp: Fast gaussian process computation
using vecchia’s approximation. 2021. Available: <https://CRAN.R-project.org/package=GpGp>
8. Pebesma E. Multivariable geostatistics in S: The gstat package. *Computers &
Geosciences*. 2004;30: 683–691.

9. Nychka D, Hammerling D, Sain S, Lenssen N. LatticeKrig: Multiresolution kriging based on markov random fields. Boulder, CO, USA: University Corporation for Atmospheric Research; 2016. doi:10.5065/D6HD7T1R
10. Venables WN, Ripley BD. Modern applied statistics with s. Fourth. New York: Springer; 2002. Available: <https://www.stats.ox.ac.uk/pub/MASS4/>
11. Finley AO, Banerjee S, Carlin BP. spBayes: An R package for univariate and multivariate hierarchical point-referenced spatial models. Journal of Statistical Software. 2007;19: 1–24. Available: <https://www.jstatsoft.org/article/view/v019i04>
12. Finley AO, Datta A, Banerjee S. spNNGP R package for nearest neighbor Gaussian process models. Journal of Statistical Software. 2022;103: 1–40. doi:10.18637/jss.v103.i05
13. Lee D. CARBayes: An R package for Bayesian spatial modeling with conditional autoregressive priors. Journal of Statistical Software. 2013;55: 1–24. Available: <https://www.jstatsoft.org/htaccess.php?volume=55&type=i&issue=13>
14. Ronnegard L, Shen X, Alam M. Hglm: A package for fitting hierarchical generalized linear models. The R Journal. 2010;2: 20–28. Available: https://journal.r-project.org/archive/2010-2/RJournal_2010-2_Roennegaard-et-al.pdf
15. Wickham H. ggplot2: Elegant graphics for data analysis. Springer-Verlag New York; 2016. Available: <https://ggplot2.tidyverse.org>
16. Appelhans T, Detsch F, Reudenbach C, Woellauer S. Mapview: Interactive viewing of spatial data in r. 2022. Available: <https://CRAN.R-project.org/package=mapview>
17. Tobler WR. A computer movie simulating urban growth in the detroit region. Economic geography. 1970;46: 234–240.
18. Anselin L, Syabri I, Kho Y. GeoDa: An introduction to spatial data analysis. Handbook of applied spatial analysis. Springer; 2010. pp. 73–89.
19. Ver Hoef JM, Peterson EE, Hooten MB, Hanks EM, Fortin M-J. Spatial autoregressive models for statistical inference from ecological data. Ecological Monographs. 2018;88: 36–59.
20. Pebesma E. Simple Features for R: Standardized Support for Spatial Vector Data. The R Journal. 2018;10: 439–446. doi:10.32614/RJ-2018-009
21. Patterson D, Thompson R. Recovery of inter-block information when block sizes are unequal. Biometrika. 1971;58: 545–554.
- 22.

- Harville DA. Maximum likelihood approaches to variance component estimation and to related problems. *Journal of the American Statistical Association*. 1977;72: 320–338.
23. Wolfinger R, Tobias R, Sall J. Computing gaussian likelihoods and their derivatives for general linear mixed models. *SIAM Journal on Scientific Computing*. 1994;15: 1294–1310.
24. Cressie N. Fitting variogram models by weighted least squares. *Journal of the international Association for mathematical Geology*. 1985;17: 563–586.
25. Curriero FC, Lele S. A composite likelihood approach to semivariogram estimation. *Journal of Agricultural, biological, and Environmental statistics*. 1999; 9–28.
26. Hoeting JA, Davis RA, Merton AA, Thompson SE. Model selection for geostatistical models. *Ecological Applications*. 2006;16: 87–98.
27. Pinheiro J, Bates D. Mixed-effects models in s and s-PLUS. Springer science & business media; 2006.
28. Hastie T, Tibshirani R, Friedman JH, Friedman JH. The elements of statistical learning: Data mining, inference, and prediction. Springer; 2009.
29. Montgomery DC, Peck EA, Vining GG. Introduction to linear regression analysis. John Wiley & Sons; 2021.
30. Myers RH, Montgomery DC, Vining GG, Robinson TJ. Generalized linear models: With applications in engineering and the sciences. John Wiley & Sons; 2012.
31. Cook RD, Weisberg S. Residuals and influence in regression. New York: Chapman; Hall; 1982.
32. Robinson D, Hayes A, Couch S. Broom: Convert statistical objects into tidy tibbles. 2021. Available: <https://CRAN.R-project.org/package=broom>
33. Box GE, Cox DR. An analysis of transformations. *Journal of the Royal Statistical Society: Series B (Methodological)*. 1964;26: 211–243.
34. Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*. 2015;67: 1–48. doi:10.18637/jss.v067.i01
35. MacQueen J, others. Some methods for classification and analysis of multivariate observations. *Proceedings of the fifth berkeley symposium on mathematical statistics and probability*. Oakland, CA, USA; 1967. pp. 281–297.