

# Analysis of univariate point referenced spatial data

---

Abhi Datta

February 1, 2018

Department of Biostatistics, Bloomberg School of Public Health, Johns Hopkins University, Baltimore, Maryland

# Review of last lecture

- Types of spatial data – point referenced, areal, point pattern
- Exploratory data analysis with point referenced data
  - Surface plots of the response, covariates and residuals
  - Empirical variograms of the residuals
- When purely covariate based models does not suffice, one needs to leverage the information from locations
  - Simple choices like adding the co-ordinates as covariates in a linear regression
  - More general model:  $y(s) = x(s)' \beta + w(s) + \epsilon(s)$  for all  $s \in D$
- How to choose the function  $w(\cdot)$ ?
- Since we want to predict at any location over the entire domain  $D$ , this choice will amount to choosing a surface  $w(s)$
- We will do this using Gaussian Processes

# Gaussian Processes (GPs)

- The collection of random variables  $\{w(s) \mid s \in D\}$  is a GP if
  - it is a **valid** stochastic process
  - all finite dimensional densities  $\{w(s_1), \dots, w(s_n)\}$  follow multivariate Gaussian distribution
- Why GPs are attractive - only need a mean function  $m(s)$  and a valid covariance function  $C(\cdot, \cdot)$
- **Advantage:** Likelihood based inference.  
 $w = (w(s_1), \dots, w(s_n))' \sim N(m, C)$  where  
 $m = (m(s_1), \dots, m(s_n))'$  and  $C = (C(s_i, s_j))$
- For the model  $y(s) = x(s)'\beta + w(s) + \epsilon(s)$ ,  $x(s)'\beta$  is **modeling the mean**. Hence,  $m(s)$  is often chosen to be 0.

## Valid covariance functions and isotropy

- $C(\cdot, \cdot)$  needs to be a **positive definite** function
- Simplifying assumptions:
  - **Stationarity**:  $C(s_1, s_2) = \text{Cov}(w(s_1), w(s_2))$  only depends on  $h = s_1 - s_2$  (and is denoted by  $C(h)$ )
  - **Isotropic**:  $C(h) = C(\|h\|)$  (**Simplest and most interpretable**)
  - **Anisotropic**: Stationary but not isotropic
- **Exponential** covariance function:  $C(h) = \sigma^2 \exp(-\phi\|h\|)$  is a **popular** choice for  $C(\cdot, \cdot)$

# Experimental evidence of BEC

- **Recall:** Empirical semivariogram:

$$\gamma(t_k) = \frac{1}{2|N(t_k)|} \sum_{s_i, s_j \in N(t_k)} (Y(s_i) - Y(s_j))^2$$

- For any stationary GP,

$$E(Y(s+h) - Y(s))^2/2 = C(0) - C(h) = \gamma(h)$$

- $\gamma(h)$  is the **semivariogram** corresponding to the covariance function  $C(h)$

- **Example:** For exponential GP,  $\gamma(t) = \begin{cases} \tau^2 + \sigma^2(1 - \exp(-\phi t)) & \text{if } t > 0 \\ 0 & \text{if } t = 0 \end{cases}$ ,  
where  $t = ||h||$

Test text for the second column.

## Notes on exponential model

$$C(t) = \begin{cases} \tau^2 + \sigma^2 & \text{if } t = 0 \\ \sigma^2 \exp(-\phi t) & \text{if } t > 0 \end{cases}.$$

- We define the **effective range**,  $t_0$ , as the distance at which this correlation has dropped to only 0.05. Setting  $\exp(-\phi t_0)$  equal to this value we obtain  $t_0 \approx 3/\phi$ , since  $\log(0.05) \approx -3$ .
- The **nugget**  $\tau^2$  is often viewed as a “**nonspatial effect variance**,”
- The **partial sill** ( $\sigma^2$ ) is viewed as a “**spatial effect variance**.”
- $\sigma^2 + \tau^2$  gives the maximum total variance often referred to as the **sill**
- Note **discontinuity** at 0 due to the nugget. **Intentional!** To account for measurement error or micro-scale variability.

# Covariance functions and semivariograms

- **Recall:** Empirical semivariogram:

$$\gamma(t_k) = \frac{1}{2|N(t_k)|} \sum_{s_i, s_j \in N(t_k)} (Y(s_i) - Y(s_j))^2$$

- For any stationary GP,

$$E(Y(s+h) - Y(s))^2/2 = C(0) - C(h) = \gamma(h)$$

- $\gamma(h)$  is the **semivariogram** corresponding to the covariance function  $C(h)$
- **Example:** For exponential GP,

$$\gamma(t) = \begin{cases} \tau^2 + \sigma^2(1 - \exp(-\phi t)) & \text{if } t > 0 \\ 0 & \text{if } t = 0 \end{cases}, \text{ where } t = ||h||$$





# The Matèrn covariance function

- The Matèrn is a very versatile family:

$$C(t) = \begin{cases} \frac{\sigma^2}{2^{\nu-1}\Gamma(\nu)} (2\sqrt{\nu}t\phi)^\nu K_\nu(2\sqrt{\nu}t\phi) & \text{if } t > 0 \\ \tau^2 + \sigma^2 & \text{if } t = 0 \end{cases}$$

$K_\nu$  is the modified Bessel function of order  $\nu$  (computationally tractable)

- $\nu$  is a smoothness parameter controlling process smoothness.  
[Remarkable!](#)
- $\nu = 1/2$  gives the exponential covariance function

## Kriging: Spatial prediction at new locations

- **Goal:** Given observations  $w = (w(s_1), w(s_2), \dots, w(s_n))'$ , predict  $w(s_0)$  for a new location  $s_0$
- If  $w(s)$  is modeled as a GP, then  $(w(s_0), w(s_1), \dots, w(s_n))'$  jointly follow multivariate normal distribution
- $w(s_0) | w$  follows a normal distribution with
  - Mean (**kriging estimator**):  $m(s_0) + c' C^{-1}(w - m)$
  - where  $m = E(w)$ ,  $C = \text{Cov}(w)$ ,  $c = \text{Cov}(w, w(s_0))$
  - Variance:  $C(s_0, s_0) - c' C^{-1} c$
- The GP formulation gives the **full predictive distribution** of  $w(s_0) | w$

## Spatial linear model

$$y(s) = x(s)' \beta + w(s) + \epsilon(s)$$

- $w(s)$  modeled as  $GP(0, C(\cdot | \theta))$  (usually without a nugget)
- $\epsilon(s) \stackrel{\text{iid}}{\sim} N(0, \tau^2)$  contributes to the nugget
- Under isotropy:  $C(s + h, s) = \sigma^2 R(\|h\| ; \phi)$
- $w = (w(s_1), \dots, w(s_n))' \sim N(0, \sigma^2 R(\phi))$  where  $R(\phi) = \sigma^2 (R(\|s_i - s_j\| ; \phi))$
- $y = (y(s_1), \dots, y(s_n))' \sim N(X\beta, \sigma^2 R(\phi) + \tau^2 I)$

## Parameter estimation

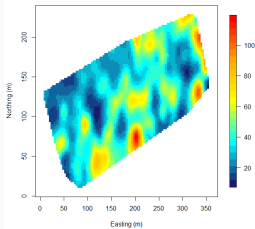
- $y = (y(s_1), \dots, y(s_n))' \sim N(X\beta, \sigma^2 R(\phi) + \tau^2 I)$
- We can obtain MLEs of parameters  $\beta, \tau^2, \sigma^2, \phi$  based on the above model and use the estimates to kriging at new locations
- In practice, the likelihood is often very **flat** with respect to the spatial covariance parameters and choice of **initial values** is important
- Initial values can be eyeballed from empirical semivariogram of the residuals from ordinary linear regression
- Estimated parameter values can be used for kriging

# Model comparison

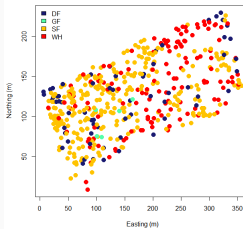
- For  $k$  total parameters and sample size  $n$ :
  - **AIC**:  $2k - 2 \log(l(y | \hat{\beta}, \hat{\theta}, \hat{\tau}^2))$
  - **BIC**:  $\log(n)k - 2 \log(l(y | \hat{\beta}, \hat{\theta}, \hat{\tau}^2))$
- Prediction based approaches using holdout data:
  - Root Mean Square Predictive Error (**RMSPE**):
$$\sqrt{\frac{1}{n_{out}} \sum_{i=1}^{n_{out}} (y_i - \hat{y}_i)^2}$$
  - Coverage probability (**CP**):  $\frac{1}{n_{out}} \sum_{i=1}^{n_{out}} I(y_i \in (\hat{y}_{i,0.025}, \hat{y}_{i,0.975}))$
  - Width of 95% confidence interval (**CIW**):
$$\frac{1}{n_{out}} \sum_{i=1}^{n_{out}} (\hat{y}_{i,0.975} - \hat{y}_{i,0.025})$$
  - The last two approaches compares the distribution of  $y_i$  instead of comparing just their point predictions

# Western Experimental Forestry (WEF) data

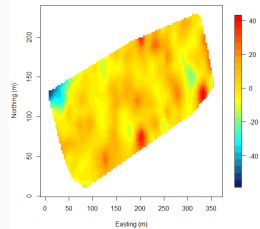
- Data consist of a census of all trees in a 10 ha. stand in Oregon
- Response of interest: Diameter at breast height (DBH)
- Covariate: Tree species (Categorical variable)



DBH



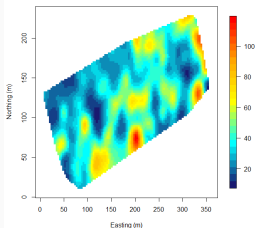
Species



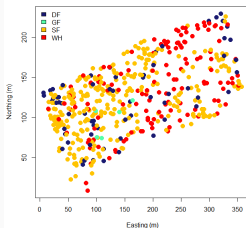
Residuals

# Western Experimental Forestry (WEF) data

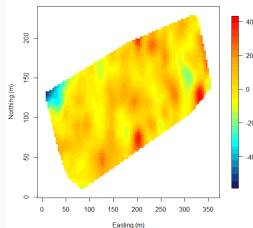
- Data consist of a census of all trees in a 10 ha. stand in Oregon
- Response of interest: Diameter at breast height (DBH)
- Covariate: Tree species (Categorical variable)



DBH



Species

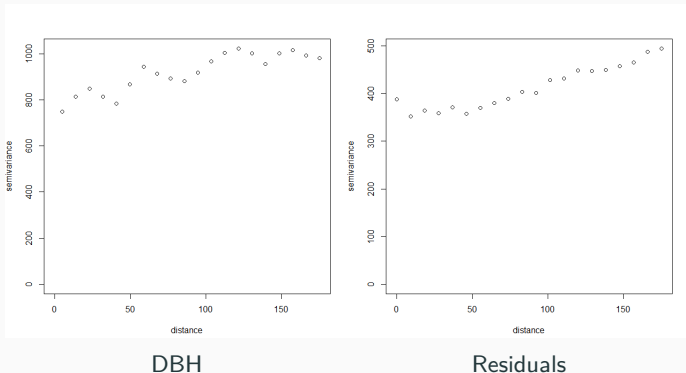


Residuals

- Local spatial patterns in the residual plot
- Simple regression on species seems to be not sufficient

# Empirical semivariograms

- Regression model:  $\text{DBH} \sim \text{Species}$



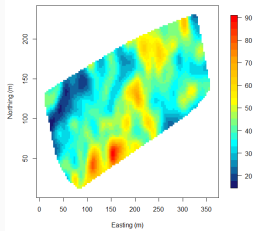
- Semivariogram of the residuals confirm **unexplained spatial variation**



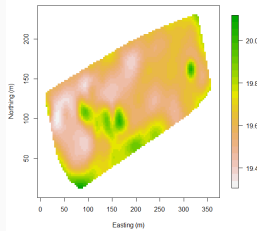
**Table:** Model comparison

	Spatial	Non-spatial
AIC	4419	4465
BIC	4448	4486
RMSPE	18	21
CP	93	93
CIW	77	82

# WEF data: Kriged surfaces



DBH Estimates



Standard errors

# Summary

- Geostatistics – Analysis of point-referenced spatial data
- Surface plots of data and residuals
- EDA with empirical semivariograms
- Modeling unknown surfaces with Gaussian Processes
- Kriging: Predictions at new locations
- Spatial linear regression using Gaussian Processes