

# Introduction to the Spatial Relationships Between Two Variables

Ronny Vallejos

Departamento de Matemática  
Universidad Técnica Federico Santa María  
Valparaíso, Chile  
[ronny.vallejos@usm.cl](mailto:ronny.vallejos@usm.cl)



Fundación Konrad Lorenz, Bogotá, Colombia  
November, 2018

Spatial Association Between Two Processes



UNIVERSIDAD TÉCNICA  
FEDERICO SANTA MARÍA

## Outline

### □ Day 1

1. Introduction and Motivation.
2. Types of Problems to be Consider in This Shortcourse
3. A Simple Solution
4. Preliminaries and Notation (Spatial Statistics context)

### □ Day 2

1. Introduction to R
2. Uploading Images in R
3. R Examples in a Spatial Context
4. The Cross-Variogram and its Properties

### □ Day 3

4. The modified  $t$  test
5. Correlation Between one Process and Several Others
6. The Effective Sample Size

# Outline

## □ Day 4

1. The Codispersion Coefficient
2. Properties and Asymptotic Results
3. Hypothesis Testing
4. Applications

## □ Day 5

1. The Effective Sample Size Revisited
2. Properties and Asymptotic Results
3. Hypothesis Testing
4. Applications

## Material of the Course

The material contained in this short course and more can be found in:

- ▣ Vallejos, R., Osorio, F., Bevilacqua, M. (2019).  
Spatial Relationships Between Two Georeferenced Variables:  
With Applications in R.  
To appear in Springer. Number of pages  $\approx$  200.
- ▣ Book homepage <http://srb2gv.mat.utfsm.cl>
- ▣ Personal webpage <http://rvallejos.mat.utfsm.cl>
- ▣ e-mail: [ronny.vallejos@usm.cl](mailto:ronny.vallejos@usm.cl)

# 1. Introduction and Motivation

- Consider two spatial processes  $\{X(\mathbf{s}) : \mathbf{s} \in D\}$  and  $\{Y(\mathbf{s}) : \mathbf{s} \in D\}$ , where  $D \subset \mathbb{R}^2$ .
- The available information is the observations  $X(\mathbf{s}_1), \dots, X(\mathbf{s}_n)$  and  $Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n)$ .
- How to quantify the linear correlation between  $X(\cdot)$  and  $Y(\cdot)$  taking into account the georeferencing information?
- Is the Pearson's correlation coefficient enough in a spatial context?

## Some Examples

### Example 1: The Pinus Radiata Dataset

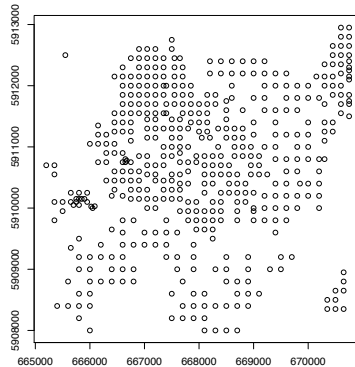
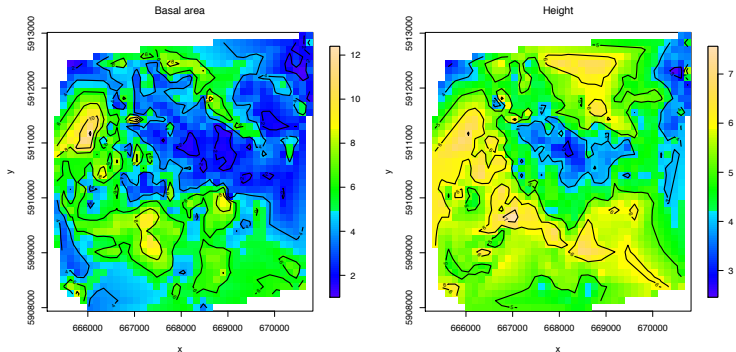


Figure: Locations where the samples were taken.

## Example 1: The Pinus Radiata Dataset



**Figure:** (Left) Bilinear interpolation of the three basal areas; (Right) Bilinear interpolation of the three heights.

## The Pinus Radiata Dataset

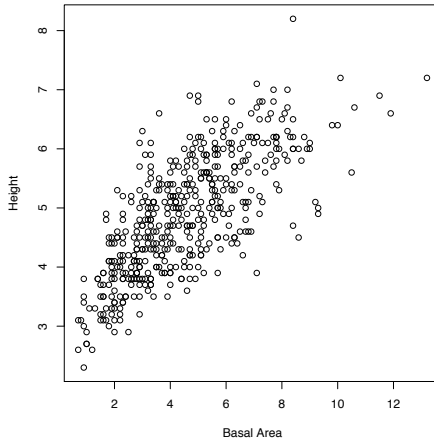
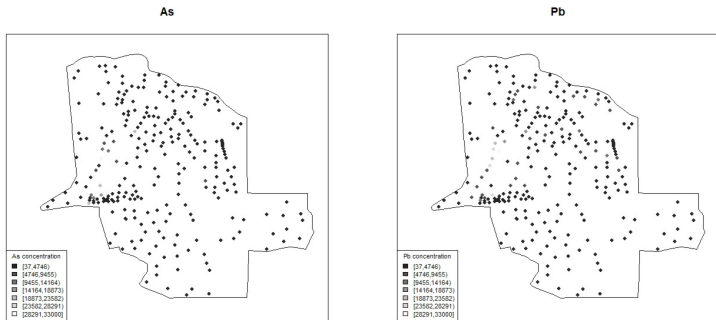


Figure: Height versus Basal Area (468 observations),  $\rho = 0.7021$ .

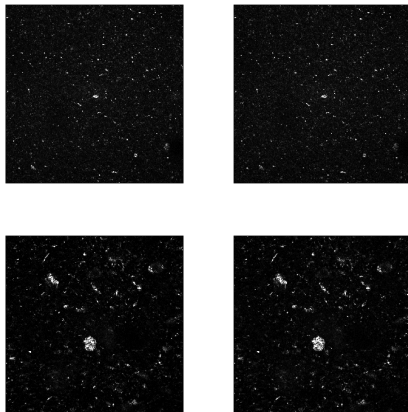


## Example 2: The Murray Smelter Site Dataset



**Figure:** Locations of 253 geocoded aggregated surface soil samples collected in a 0.5 square mile area in Murray, Utah and their measured concentrations of As and Pb,  $\rho = 0.5893$ .

## Example 3: Similarity Between Images



**Figure:** Dispersion of nanotubes. Images taken at NIST, USA. The images of the top were taken at the same distance. The images of the bottom were taken at the same distance but closer.

## Example 4: Directional Contamination



**Figure:** (Left) Original image (Lenna); (Right) Image transformed into the direction  $h = (1, 1)$ .

## Philosophy of the Course

Throughout, I have attempted to follow one basic principle:  
never give an estimator without giving a confidence set.

Larry Wasserman  
Department of Statistics  
Carnegie Mellon University



## Easy Solution

### The $t$ Test

- Arrange the observations of each variable in a column. Under normality assume that  $\text{cor}[X(s), Y(s)] = \rho$  and consider the hypothesis testing problem:

$$H_0 : \rho = 0 \text{ versus } H_1 : \rho \neq 0.$$

- Use the statistic

$$t = r \frac{\sqrt{n-2}}{\sqrt{1-r^2}},$$

where  $r$  is an estimator of  $\rho$ .

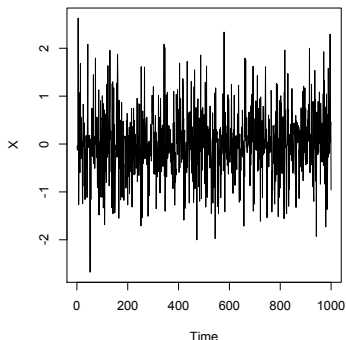
- Compute the  $p$ -value associated with the test.
- For a fixed level  $\alpha$  we reject  $H_0$  if  $p < \alpha$ , otherwise  $H_0$  is not rejected.

## 2. Preliminaries and Notation

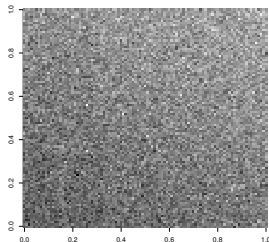
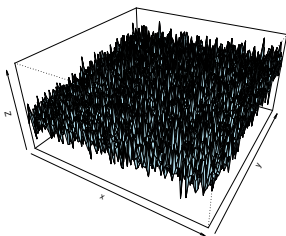
### Definition

Let  $(\Omega, \mathcal{F}, \mathcal{P})$  be a probability space, and let  $D$  be an arbitrary index set. A spatial process is a function  $X : (\Omega, \mathcal{F}, \mathcal{P}) \times D \rightarrow \mathbb{R}$ , such that for all  $s \in D$ ,  $X(s)$  is a random variable. Commonly  $D \subset \mathbb{R}^d$ .

$X(s) = A \cos(\eta s + \phi)$ , where  $A$  is a random variable independent of  $\phi \sim \mathcal{U}(0, 2\pi)$ , and  $\eta$  is a fixed constant. For the particular case when  $A \sim \mathcal{N}(0, 1)$  and  $\eta = 1$ , 1000 observations were generated.



## Realizations of a Spatial Process



**Figure:** A realization from a spatial process defined on a finite grid of  $D = \mathbb{Z}^2$ .  $Z(x, y) = \beta_1 x + \beta_2 y + \epsilon(x, y)$ , where  $\{\epsilon(x, y)\}$  is a collection of independent and identically distributed random variables with zero mean and variance  $\sigma^2$ .

## Stationary Processes

- A condition that guaranties the existence of the mean, variance and covariance functions of a spatial process is the second order property:

$$\mathbb{E}[X^2(\mathbf{s})] < \infty, \text{ for all } \mathbf{s} \in D.$$

### Theorem

Let  $X$  be a second order spatial process . Then

- ▶  $\mathbb{E}[X(\mathbf{s})] < \infty$ , para todo  $\mathbf{s} \in D$ .
  - ▶  $\text{var}[X(\mathbf{s})] < \infty$ , para todo  $\mathbf{s} \in D$ .
  - ▶  $\text{cov}[X(\mathbf{s}), X(\mathbf{t})] < \infty$ , para todo  $\mathbf{s}, \mathbf{t} \in D$ .
- A second order spatial process is said weakly stationary if
- i)  $\mathbb{E}[X(\mathbf{s})] = \mu$ , for all  $\mathbf{s} \in D$ .
  - ii)  $\text{cov}[X(\mathbf{s}_i), X(\mathbf{s}_j)] = g(\mathbf{s}_i - \mathbf{s}_j) = C(\mathbf{h})$ , for all  $\mathbf{s}_i, \mathbf{s}_j \in D$  and for some function  $g$ . An immediate consequence is that  $\text{var}[X(\mathbf{s})] = g(\mathbf{0})$ , constant with respect to  $\mathbf{s}$ .



## Covariance Functions

The covariance function of a stationary process is denoted as

$$C(\mathbf{h}) = \text{cov}[X(\mathbf{s}), X(\mathbf{s} + \mathbf{h})], \quad \mathbf{s}, \mathbf{h} \in D.$$

- The covariance function of a weakly stationary process must satisfy the positive-definiteness condition, i.e.,

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j C(\mathbf{s}_i - \mathbf{s}_j) \geq 0,$$

for all  $\mathbf{s}_i, \mathbf{s}_j \in D$  and for all  $a_i, a_j \in \mathbb{R}$ .

- The corresponding covariance matrix of the vector  $(X(\mathbf{s}_1), \dots, X(\mathbf{s}_n))^T$  denoted as  $\Sigma$ , satisfies that its  $ij$ -th element is

$$\Sigma_{ij} = C(\mathbf{s}_i - \mathbf{s}_j).$$

## Covariance Functions

### Proposition

Let  $C$  be the covariance function of a stationary process. Then

- $C(\mathbf{h}) = C(-\mathbf{h}), \forall \mathbf{h} \in D$ . (even function)
- $|C(\mathbf{h})| \leq C(\mathbf{0}), \forall \mathbf{h} \in D$ . (bounded function)
- If  $C$  is a valid covariance function in  $\mathbb{R}^d$ , then  $C$  is valid in  $\mathbb{R}^k$  for all integer  $k < d$ .
- If  $C_j(\mathbf{h})$  are valid covariance functions,  $j = 1, \dots, k$ , then  $\sum_{j=1}^k b_j C_j(\mathbf{h})$  is a valid covariance function, if  $b_j \geq 0$ , for all  $j$ .
- If  $C_j(\mathbf{h})$  are valid covariance functions,  $j = 1, \dots, k$ , then  $\prod_{j=1}^k C_j(\mathbf{h})$  is also a valid covariance function.

## Strongly Stationary Processes

### Definition

A spatial process  $X$  is said strictly stationary if  $\forall k \in \mathbb{N}$ ,  $\forall (s_1, \dots, s_k) \in D^k$ , and  $\forall \mathbf{h} \in D$ , the distribution of the vector  $(X(s_1 + \mathbf{h}), \dots, X(s_k + \mathbf{h}))^\top$  is independent of  $\mathbf{h}$ .

### Theorem

*Let  $X$  be a second order process.*

*$X$  strictly stationary  $\Rightarrow X$  weakly stationary*

### Definition

A spatial process  $X$  is Gaussian if,  $\forall k \in \mathbb{N}$ , the vector  $(X(s_1), \dots, X(s_k))^\top$  has a multivariate normal distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ . We will use the notation  $(X(s_1), \dots, X(s_k))^\top \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

### Theorem

*$X$  strictly stationary  $\iff X$  weakly stationary and Gaussian*

## Intrinsically Stationary Processes

### Definition

A spatial process  $\{X(\mathbf{s}) : \mathbf{s} \in D\}$  is called intrinsically stationary if for all  $\mathbf{s} \in D$ ,  $\mathbb{E}[X(\mathbf{s})] = \mu$ , and

$$\text{var}[X(\mathbf{s} + \mathbf{h}) - X(\mathbf{s})] = 2\gamma(\mathbf{h})$$

is only a function of  $\mathbf{h}$ . In such a case, the function  $2\gamma(\mathbf{h})$  is called a variogram of the process and  $\gamma(\mathbf{h})$  is called a semivariogram.

- If  $\mathbb{E}[X(\mathbf{s} + \mathbf{h}) - X(\mathbf{s})] = 0$ , then
$$\text{var}[X(\mathbf{s} + \mathbf{h}) - X(\mathbf{s})] = 2\gamma(\mathbf{h}) = \mathbb{E}[X(\mathbf{s} + \mathbf{h}) - X(\mathbf{s})]^2.$$
- In practice, small variance is expected for points that are close in space, while large variance is expected for points with a large separation in space according to the **first law of geography** (Tobler, 1970).

## Properties of the Semi-Variogram

The semi-variogram  $\gamma$  of an intrinsically stationary process  $X$  satisfies:

- $\gamma(\mathbf{h}) = \gamma(-\mathbf{h})$  (even function) and  $\gamma(\mathbf{0}) = 0$ ;
- If  $A$  is a linear function defined on  $\mathbb{R}^d$ , then the function  $h \rightarrow \gamma(A\mathbf{h})$  it is also a semi-variogram.
- If  $\gamma$  is continuous in  $\mathbf{0}$ , then  $\gamma$  is continuous for all  $\mathbf{h}$   $\gamma$  is locally bounded.
- If  $\gamma$  is bounded in a neighborhood of  $\mathbf{0}$ , then there exist  $a, b \in \mathbb{R}$  such that for all  $x \in D$ ,  $\gamma(x) \leq a\|x\|^2 + b$ .
- **Conditionally Negative Definite.** For all  $n \in \mathbb{N}$ , for all  $\mathbf{a} \in \mathbb{R}^n$  such that  $\sum_{i=1}^n a_i = 0$ , and for all  $(s_1, \dots, s_n) \in D^n$ , 
$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j \gamma(s_i - s_j) \leq 0.$$
- A continuous function  $\gamma$  on  $\mathbb{R}^d$  such that  $\gamma(\mathbf{0}) = 0$  is a semi-variogram iff for all  $a > 0$ ,  $h \rightarrow e^{-a\gamma(h)}$  is a covariance function.

## Covariance and Semi-Variograms

- For weakly stationary processes, there is a relationship between  $\gamma(\mathbf{h})$  and  $C(\mathbf{h})$ :

$$\begin{aligned}2\gamma(\mathbf{h}) &= \text{var}[X(\mathbf{s} + \mathbf{h}) - X(\mathbf{s})] \\&= \text{var}[X(\mathbf{s} + \mathbf{h})] + \text{var}[X(\mathbf{s})] - 2\text{cov}[X(\mathbf{s} + \mathbf{h}), X(\mathbf{s})] \\&= 2(C(\mathbf{0}) - C(\mathbf{h})).\end{aligned}$$

Thus,

$$\gamma(\mathbf{h}) = C(\mathbf{0}) - C(\mathbf{h}).$$

### Example

- Consider the AR(1) process  $X(t) = \phi X(t-1) + \epsilon(t)$ , where  $|\phi| < 1$ , and  $\epsilon(t) \sim WN(0, \sigma^2)$ .
- $C(h) = \sigma^2 \phi^h / (1 - \phi^2)$ ,  $C(0) = \sigma^2 / (1 - \phi^2)$ .
- $\gamma(h) = \sigma^2 (1 - \phi^h) / (1 - \phi^2)$ .

## Covariance and Semi-Variogram Plots

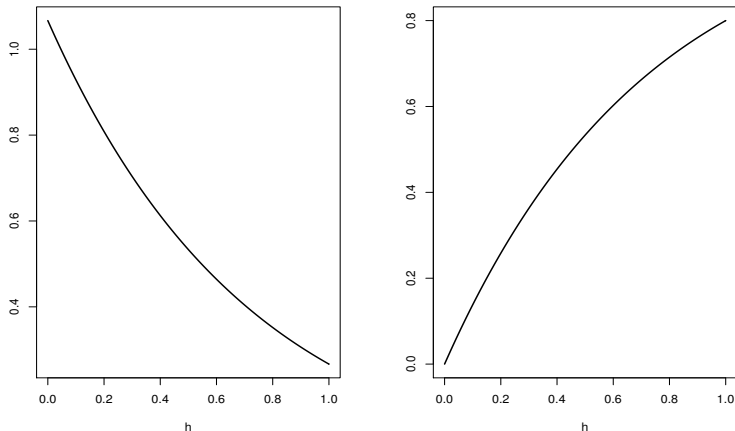


Figure:  $C(h)$  and  $\gamma(h)$  for an AR(1) model.

## Isotropy

- If a variogram (semivariogram) depends on  $\|h\|$ , the process  $\{X(s)\}$  is called isotropic; otherwise, the process  $\{X(s)\}$  is called anisotropic.

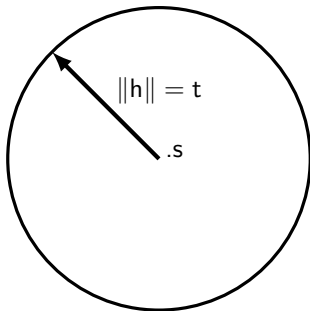


Figure: Circular (isotropic) correlation structure.  
Spatial Association Between Two Processes





## Typical Semi-Variogram Plot

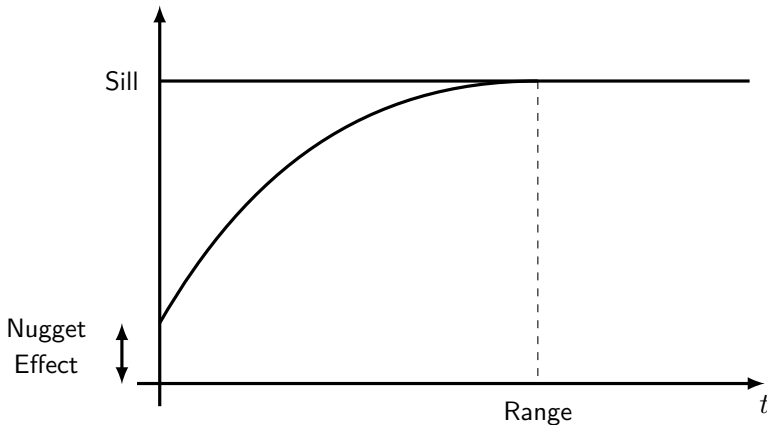


Figure: Behavior of a typical semivariogram model.  $\|h\| = t$

## Parametric Covariance Models

- There are several parametric covariance and semi-variogram models that are valid.
- The Matérn covariance family is described by

$$C(t) = \frac{\sigma^2}{2^{\nu-1}\Gamma(\nu)} \left(\frac{t}{r}\right)^{\nu} \kappa_{\nu} \left(\frac{t}{r}\right),$$

where  $\kappa_{\nu}$  is a modified Bessel function of the second kind of order  $\nu$ ,  $\Gamma$  is the gamma function,  $r > 0$  is the range, and  $\nu > 0$  is the smoothness parameter.

- If  $\nu = 1/2$ , then  $C(t) = \sigma^2 \exp(-t/r)$ .
- The covariance matrix of a spatial process with Matérn covariance function is positive definite.

# Parametric Covariance Models

Model	Covariance Function $C(t)$
Lineal	$C(t)$ does not exist.
Spheric	$C(t) = \begin{cases} 0, & \text{if } t \geq 1/\phi, \\ \sigma^2 \left( 1 - \frac{3}{2}\phi t + \frac{1}{2}(\phi t)^3 \right), & \text{if } 0 < t < 1/\phi, \\ \tau^2 + \sigma^2, & \text{otherwise.} \end{cases}$
Exponential	$C(t) = \begin{cases} \sigma^2 \exp(-\phi t), & \text{if } t > 0, \\ \tau^2 + \sigma^2, & \text{otherwise.} \end{cases}$
Gaussian	$C(t) = \begin{cases} \sigma^2 \exp(-(\phi t)^2), & \text{if } t > 0, \\ \tau^2 + \sigma^2, & \text{otherwise.} \end{cases}$
Wave	$C(t) = \begin{cases} \sigma^2 \frac{\sin(\phi t)}{\phi t}, & \text{if } t > 0, \\ \tau^2 + \sigma^2, & \text{otherwise.} \end{cases}$
Matérn	$C(t) = \begin{cases} \frac{\sigma^2}{2^{\nu-1}\Gamma(\nu)} (2\sqrt{\nu}t\phi)^{\nu} K_{\nu}(2\sqrt{\nu}t\phi), & \text{if } t > 0, \\ \tau^2 + \sigma^2, & \text{otherwise.} \end{cases}$

## Estimation of the Variogram

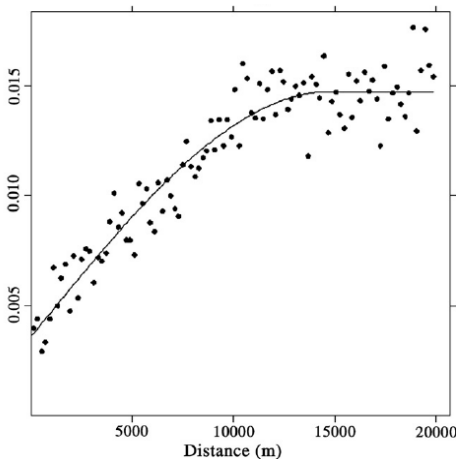
- For  $n$  sampling sites  $\mathbf{s}_1, \dots, \mathbf{s}_n$ , a natural and unbiased estimator of the semivariogram of a spatial process  $\{X(\mathbf{s}) : \mathbf{s} \in D \subset \mathbb{R}^2\}$  is

$$\hat{\gamma}(\mathbf{h}) = \frac{1}{2|N(\mathbf{h})|} \sum_{N(\mathbf{h})} (X(\mathbf{s}_i) - X(\mathbf{s}_j))^2,$$

where  $N(\mathbf{h}) = \{(\mathbf{s}_i, \mathbf{s}_j) : \mathbf{s}_i - \mathbf{s}_j = \mathbf{h}, 1 \leq i, j \leq n\}$ , and  $|\cdot|$  denotes the cardinality of a set.

- Robust estimators have been studied by Cressie (1980), Genton (1998) and García-Soidán (2004) among others.
- the computation of the semivariogram is available in the R packages GeoR, RandomFields, and sgeostat.

## Estimation of the Variogram



Spatial Association Between Two Processes



## Estimation

- Let us assume a parametric model for the semi-variogram of the form  $\{\gamma(\cdot, \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p\}$ . Commonly  $\boldsymbol{\theta} = (\sigma^2, \tau, \phi)$ .
- OLS estimation.** The least squares estimator of  $\boldsymbol{\theta}$  is:

$$\hat{\boldsymbol{\theta}}_{OLS} = \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} \sum_{k=1}^K (\hat{\gamma}(\mathbf{h}_k) - \gamma(\mathbf{h}_k))^2$$

- GLS Estimator.** The generalized least squares estimator of  $\boldsymbol{\theta}$  is:

$$\boldsymbol{\theta}_{GLS} = \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} (\hat{\gamma}_K - \gamma_K(\boldsymbol{\theta}))^\top \Sigma_{\hat{\gamma}_K}^{-1}(\boldsymbol{\theta}) (\hat{\gamma}_K - \gamma_K(\boldsymbol{\theta})),$$

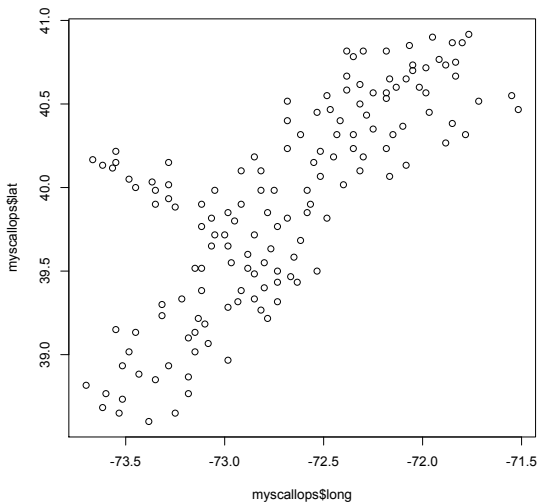
where  $\hat{\gamma}_K = (\hat{\gamma}(\mathbf{h}_1), \dots, \hat{\gamma}(\mathbf{h}_K))^\top$ ,  
 $\gamma_K(\boldsymbol{\theta}) = (\gamma(\mathbf{h}_1, \boldsymbol{\theta}), \dots, \gamma(\mathbf{h}_K, \boldsymbol{\theta}))^\top$  and  $\Sigma_{\hat{\gamma}_K}^{-1}(\boldsymbol{\theta})$  is the inverse of the covariance matrix of  $\hat{\gamma}_K(\boldsymbol{\theta})$ .

## An Example: Scallops Dataset

- Scallops (ostiones) is a dataset collected in the atlantic ocean in the north part of the United States.
- 148 observations (latitud, longitud, catches) were collected in total.
- The sampling procedure was carried out between 1990 and 1993.
- The dataset has been provided by the National Marine Fisheries Service of the USA.

Banerjee et al. (2004). Hierarchical Modelling and Analysis for Spatial Data. Chapman & Hall/CRC, Boca Ratón.

## An Example: Scallops



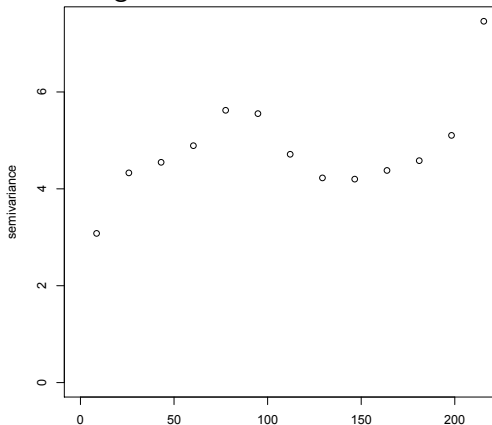
Spatial Association Between Two Processes





## An Example: Scallops

### Empirical Semi-Variogram



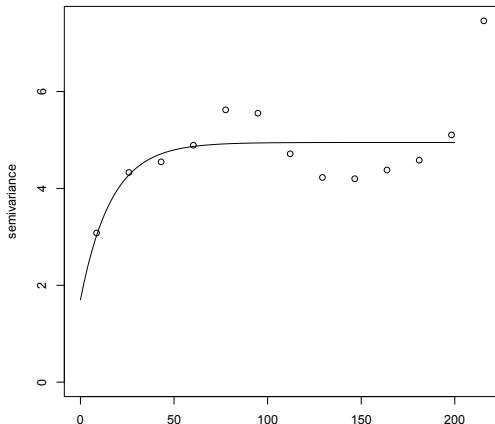
Variograma clasico

Spatial Association Between Two Processes



## An Example: Scallops

### Exponencial Semi-Variogram

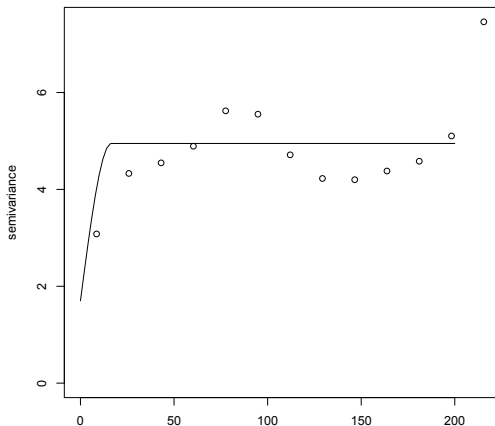


Spatial Association Between Two Processes



# An Example: Scallops

## Spherical Semi-Variogram

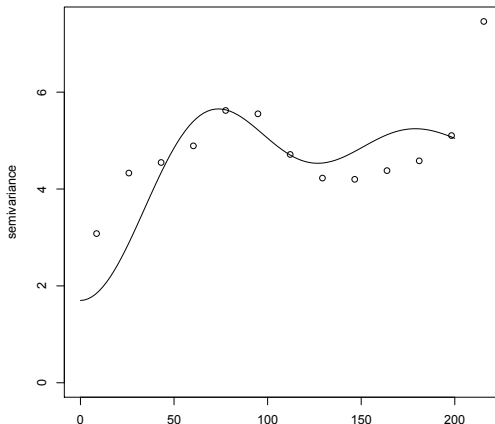


Spatial Association Between Two Processes



## Un Ejemplo: Scallops

Wave Model:  $\gamma(\mathbf{h}) = c \cdot \sin(a\|\mathbf{h}\|)/(a\|\mathbf{h}\|)$ , for  $\|\mathbf{h}\| > 0$ .

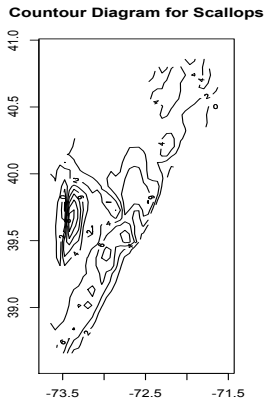
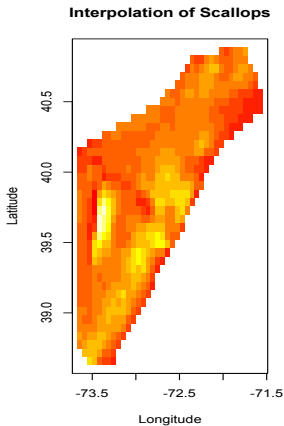


Spatial Association Between Two Processes



# Semi-Variogram Applications

## Kriging interpolation



Spatial Association Between Two Processes



# Day 2

