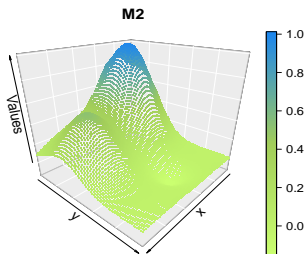


# Statistical methods for spatial data analysis

## Lecture 9: Interpolation of spatial data

Joaquin Cavieres

Geoinformatics, Bayreuth University



# Outline

1

## Introduction

- Variability
- Stochastic methods

2

## Geostatistics

- Gaussian random fields

3

## Kriging

- Estimation of  $\mu(s)$
- Ordinary Kriging (OK)

# 1. Introduction

Geostatistical data are data that could in principle be measured anywhere, but that typically come as measurements at a limited number of observation locations: think of gold grades in an ore body or particulate matter in air samples.

The pattern of observation locations is usually not of primary interest, as it often results from considerations ranging from economical and physical constraints to being 'representative' or random sampling varieties.

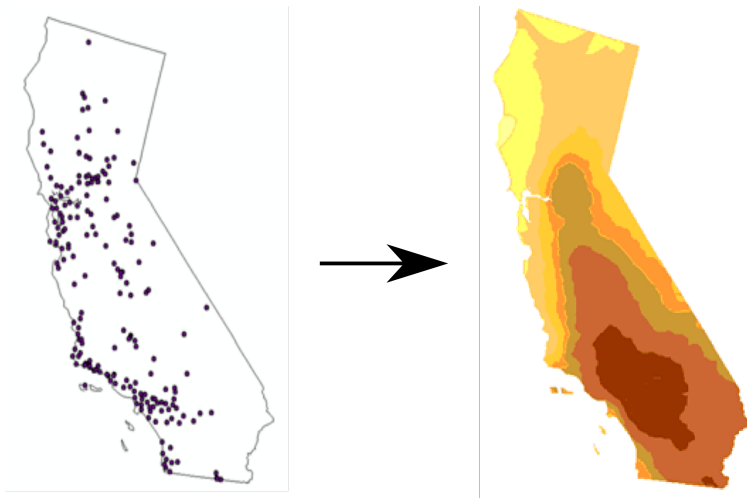
Problems what we can solve using Geostatistics are:

- The estimation of ore grades over mineable units, based on drill hole data
- Interpolation of environmental variables from sample or monitoring network data (e.g. air quality, soil pollution, ground water head, hydraulic conductivity)
- Interpolation of physical or chemical variables from sample data

But, what means interpolation?

But, what means interpolation?

Spatial interpolation is the prediction of a given phenomenon in unmeasured locations. In this case we need a spatial interpolation model to calculate predicted values of the variable of interest, given some observed data.



**Figure 1:** Spatial interpolation (Point locations of ozone monitoring stations, Interpolated prediction surface) (<https://geobgu.xyz/r-2020/spatial-interpolation-of-point-data.html>).



We can split spatial interpolation in two categories:

- **Deterministic models**: we use arbitrary parameter values, for example: IDW, NNI.
- **Statistical models**: we use parameters chosen objectively based on the data, for example: Kriging

And, what is Geostatistics?

And, what is Geostatistics? Well, first we have to consider two important concepts:

- Variability
- Stochastic methods

## 1.1. Variability

Natural phenomena generally occur in space, time or space/time. If we consider the space (spatial domain), for example, the topographic surface or a groundwater contamination one can observe high variability within small distances. However, this type of processes can not be described based on physical and chemical laws completely.

So, considering the underlying variability of the process, then we have to evaluate the following:

- Laboratory measurements are necessary to try quantify (in some way) the phenomena
- There is always a limited degree of explanation one can achieve since upscaling is non trivial

## 1.2. Stochastic methods

Usually we can use probabilistic and statistical methods for describing partly known (or sampled) natural parameters.

Measurement values of a certain parameter, obtained from different locations are often treated as different outcomes of the same random variable.

So, what is Geostatistics?

So, what is Geostatistics?

### General definition

Geostatistics can be viewed as a collection of all statistical and probabilistic methods applied in geosciences.



# What are Geostatistical data?

## What are Geostatistical data?

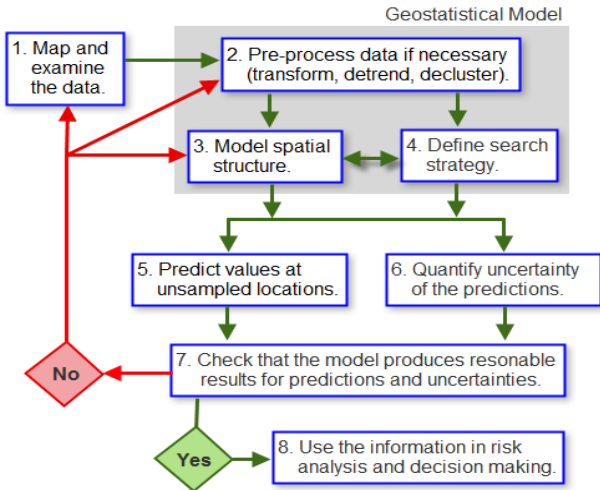
### Definition

Geostatistical data are measurements about a spatially continuous phenomenon that have been collected at particular sites.

This type of data may represent, for example, the temperature registered at several monitoring stations, the rain in some specific spatial area or the pollution of some pollutant.

The common steps to interpolate an spatial model for geostatistical data are:

- Analyse the data.
- Calculating the empirical semivariogram or covariance values.
- Fit (train) a model to the observed values.
- Predict values by the Kriging method and quantify the error (uncertainty) associated with them for every location in the output surface.



**Figure 2:** Source: <http://pro.arcgis.com/en/pro-app/help/analysis/geostatistical-analyst/the-geostatistical-workflow.htm>

For spatial interpolation, in this case we only will consider the statistical models for Geostatistical data.

## 2. Geostatistics

Let  $Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n)$  are observations of a spatial variable  $Y$  at the spatial locations  $\mathbf{s}_1, \dots, \mathbf{s}_n$ . Thus, we can assume that the spatial (geostatistical) data is a realization of a spatial stochastic process:

$$\{Y(\mathbf{s}) : \mathbf{s} \in D \subset \mathbb{R}^2\}, \quad (1)$$

where  $D$  is a fixed subset  $\mathbb{R}^2$  and the spatial index  $\mathbf{s}$  varies continuously throughout  $D$ .

## 2.1. Gaussian random fields

A Gaussian random field (GRF) is a collection of random variables where the observations occur in a continuous domain, and the collection of the random variables has a multivariate normal distribution.

$$\{Y(\mathbf{s}) : \mathbf{s} \in \mathcal{D} \subset \mathbb{R}^2\} \quad (2)$$

Additionally, a stochastic (spatial) process is strictly stationary if for a set of locations  $\mathbf{s}_i$ ,  $i = 1, \dots, n$ , it has a constant mean and invariant to translations, that is, the distribution of  $\{Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n)\}$  is the same that of  $\{Y(\mathbf{s}_1 + \mathbf{h}), \dots, Y(\mathbf{s}_n + \mathbf{h})\}$  for any  $\mathbf{h} \in \mathbb{R}^2$ .



A less restrictive condition is given by the second-order stationarity. Under this condition, the spatial process has constant mean

$$E[u(\mathbf{s})] = \mu, \forall \mathbf{s} \in \mathcal{D}, \quad (3)$$

and the covariance

$$\text{Cov}(u(\mathbf{s}), u(\mathbf{s} + \mathbf{h})) = C(\mathbf{h}), \forall \mathbf{s} \in \mathcal{D}, \forall \mathbf{h} \in \mathbb{R}^2, \quad (4)$$

only depending on the differences between locations. The covariance matrix of a GRF specifies its dependence structure and it is constructed from a covariance function.

In addition, if the covariances are functions only of the distances between locations and not of the directions, the process is called **isotropic**. Other way, it is **anisotropic**. A process is said to be intrinsically stationary if in addition to the constant mean assumption it satisfies

$$\text{Var}[Y(\mathbf{s}_i) - Y(\mathbf{s}_j)] = 2\gamma(\mathbf{s}_i - \mathbf{s}_j), \forall \mathbf{s}_i, \mathbf{s}_j. \quad (5)$$

where  $2\gamma(\cdot)$  is the **variogram** and  $\gamma(\cdot)$  is the **semivariogram** ([3]).

Considering the assumption of intrinsic stationarity, the constant-mean assumption implies:

$$2\gamma(\mathbf{h}) = \text{Var}(Y(\mathbf{s} + \mathbf{h}) - Y(\mathbf{s})) = E[(Y(\mathbf{s} + \mathbf{h}) - Y(\mathbf{s}))^2], \quad (6)$$

and the semivariogram (empirical):

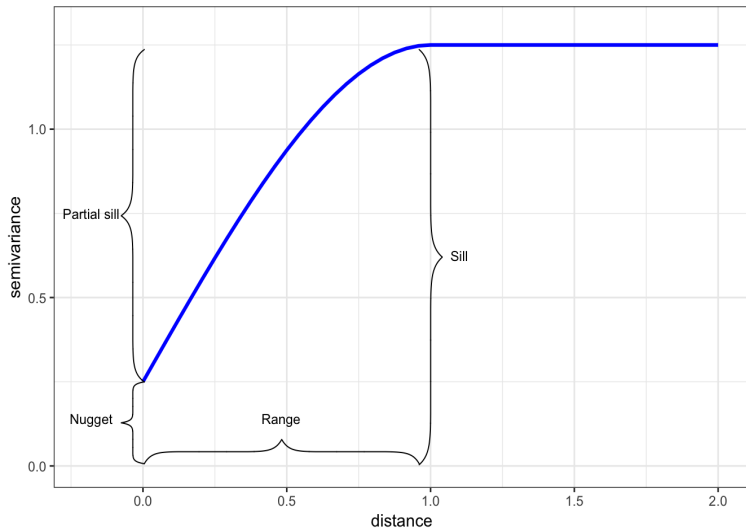
$$2\hat{\gamma}(\mathbf{h}) = \frac{1}{|N(\mathbf{h})|} \sum_{N(\mathbf{h})} (Y(\mathbf{s}_i) - Y(\mathbf{s}_j))^2, \quad (7)$$

where  $|N(\mathbf{h})|$  is the number of distinct pairs in  $N(\mathbf{h}) = \{(\mathbf{s}_i, \mathbf{s}_j) : \mathbf{s}_i - \mathbf{s}_j = \mathbf{h}, i, j = 1, \dots, n\}$ . If the process is isotropic, the semivariogram is a function of the distance  $h = \|\mathbf{h}\|$

Once calculated the variogram, it is standard practice to smooth the empirical semivariogram by fitting a parametric model to it, and the parameters used to this are: the *Nugget*, the *Sill* and the *Range*.

- Nugget: The random error process indicated by the height of the jump of the semivariogram at the discontinuity at the origin.
- Sill: The limit of the variogram tending to infinity lag distances.
- Range: The distance in which the difference of the variogram from the sill becomes negligible.

### Nugget, Sill and Range of a variogram



**Figure 3:** Nugget, Sill and Range of a variogram (Source: Freie Universitat Berlin)

There are different parametric models to fit a variogram, for example:

- Spherical:  $\gamma(h, \theta) = \begin{cases} \theta_1(\frac{3h}{2\theta_2} - \frac{h^3}{2\theta_2^3}) & \text{for } 0 \leq h \leq \theta_2 \\ \theta_1 & \text{for } h > \theta_2 \end{cases}$
- Exponential:  $\gamma(h, \theta) = \theta_1\{1 - \exp(-h/\theta_2)\}$
- Gaussian:  $\gamma(h, \theta) = \theta_1\{1 - \exp(-h^2/\theta_2^2)\}$
- Matérn:  $\gamma(h, \theta) = \theta_1 \left(1 - \frac{(h/\theta_2)^\nu \kappa_\nu(h/\theta_2)}{2^{\nu-1} \Gamma(\nu)}\right),$

where  $K_\nu(\cdot)$  is the modified Bessel function of the second kind of order  $\nu$ .

Fitting a variogram (see the example in R)

### 3. Kriging



Spatial prediction refers to the prediction of unknown quantities  $Y(s_0)$ , based on sample data  $Y(s_i)$  and assumptions regarding the form of the trend of  $Y$  and its variance and spatial correlation.

If we think about the behaviour of  $Y$ , in this quantity the large scale and small scale contribute to its variation. So, continuing with this idea, we can model this stochastic spatial process (**spatial random field**)  $\{Y(\mathbf{s}), \mathbf{s} \in D \subset \mathbb{R}^d\}$  as:

$$Y(\mathbf{s}) = \mu(\mathbf{s}) + \epsilon(\mathbf{s}), \quad \mathbf{s} \in D, \quad (8)$$

where  $\mu(\mathbf{s}) = E(Y(\mathbf{s}))$ , the mean function, and  $\epsilon(\mathbf{s})$  a zero-mean random error process.

- $\mu(\mathbf{s})$  is the first-order structure.
- $\epsilon(\mathbf{s})$  is the the second-order structure.

Here  $\mu(\mathbf{s})$  accounts for large-scale spatial variation (*global trend*), and  $\epsilon(\mathbf{s})$  for the small-scale spatial variation (*spatial dependence*). As  $\epsilon(\mathbf{s})$  typically has no spatial structure we explicitly separate it from the spatially dependent component. Thus, the geostatistical model can be written as:

$$Y(\mathbf{s}) = \mu(\mathbf{s}) + \eta(\mathbf{s}) + \epsilon(\mathbf{s}) \quad (9)$$

where

- $\eta(\mathbf{s})$  is the spatially dependent component and the spatially uncorrelated mean zero errors.
- $\epsilon(\mathbf{s})$  is the measurement error.

Note that the processes  $\eta(\mathbf{s})$  and  $\epsilon(\mathbf{s})$  are independent, and it is often referred to as a *nugget effect*.

## 3.1. Estimation of $\mu(s)$

The mean function  $\mu(s)$  is specified by a parametric model,  $\mu(s; \beta)$ , then:

$$\mu(s; \beta) = \mathbf{X}(s)^T \beta \quad (10)$$

where  $\mathbf{X}(s)$  is a vector of covariates (explanatory variables) observed at  $s$ , and  $\beta$  is a parameter vector.

The standard method for fitting the mean function  $\mu(\mathbf{s})$  is the ordinary least squares (OLS). The mathematical expression is the following:

$$\hat{\beta}_{OLS} = \operatorname{argmin} \sum_{i=1}^n [Y(\mathbf{s}_i) - \mathbf{X}(\mathbf{s}_i)^T \beta]^2 \quad (11)$$

Fitted values and fitted residuals at data locations are given by the equation  $\hat{Y} = \mathbf{X}^T \hat{\beta}_{OLS}$  and  $\hat{\epsilon} = Y - \hat{Y}$ .

## 3.2. Ordinary Kriging (OK)

We will consider the model:

$$Y(\mathbf{s}) = \mu + \epsilon(\mathbf{s}) \quad (12)$$

where  $\mu$  is the constant stationary function (global mean) and  $\epsilon(\mathbf{s})$  is the spatially correlated stochastic part of variation

The predictions are done using the following equation:

$$\hat{Y}_{OK}(\mathbf{s}_0) = \sum_{i=1}^n w_i(\mathbf{s}_0) \cdot Y(\mathbf{s}_i) = \boldsymbol{\lambda}_0^T \cdot \mathbf{Y} \quad (13)$$

where  $\boldsymbol{\lambda}_0^T$  is the vector of Kriging parameters (weights,  $(w_i)$ ),  $\mathbf{Y}$  is the vector of  $n$  observations at primary locations.

See the example in R



# Thanks!...



Moraga, P., 2023. Spatial Statistics for Data Science: Theory and Practice with R



Pebesma, E., & Bivand, R. (2023). Spatial Data Science: With Applications in R. CRC Press.



Cressie, N. (2015). Statistics for spatial data. John Wiley & Sons.



Bivand, R. S., Pebesma, E. J., Gomez-Rubio, V., & Pebesma, E. J. (2008). Applied spatial data analysis with R (Vol. 747248717, pp. 237-268). New York: Springer.