**ISRIC Spring School 2018**

# Introduction to Geostatistics
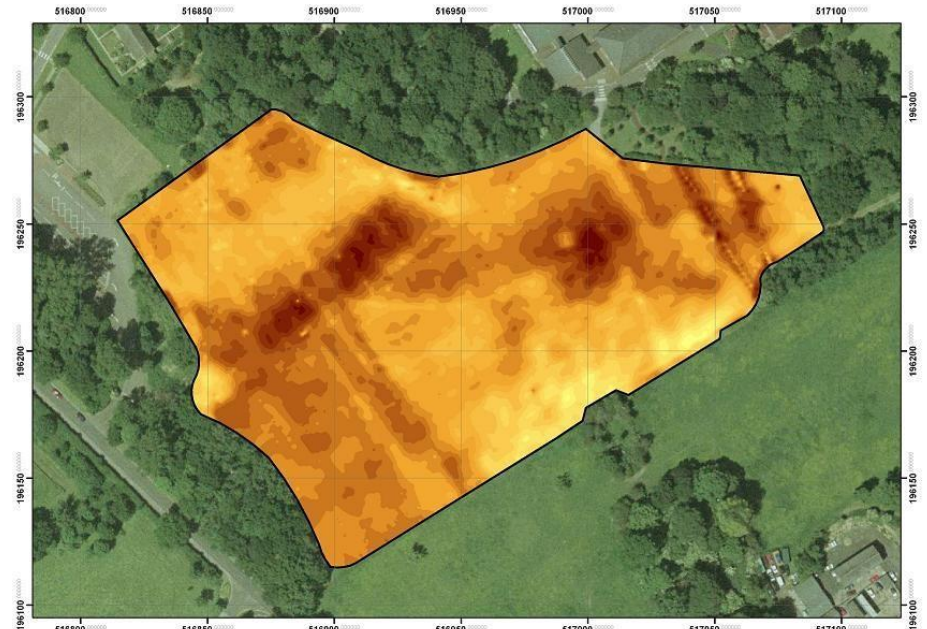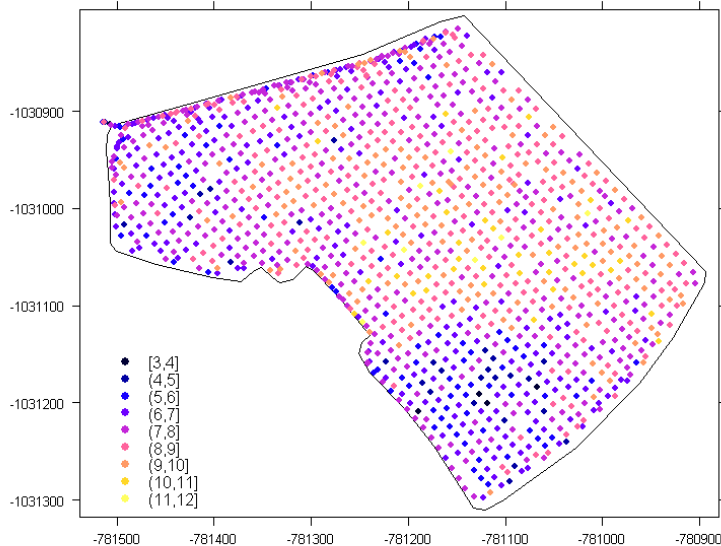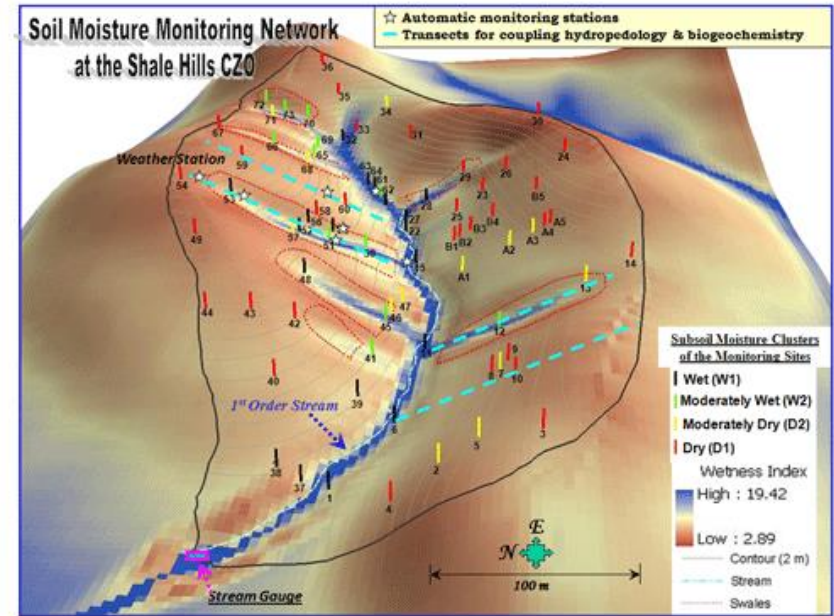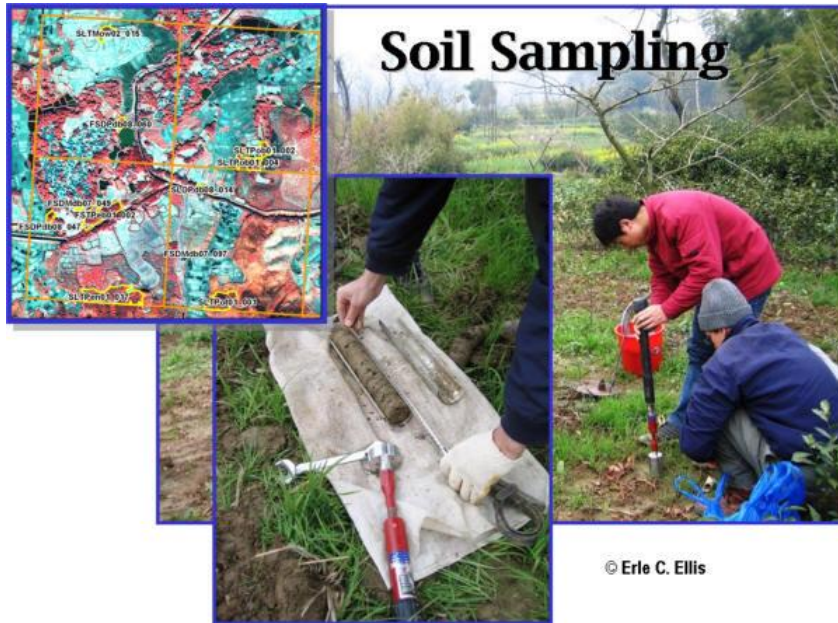
ISRIC — World Soil Information

## Bas Kempen

# Outline of this lecture

- Explore spatial variation
- Quantify spatial variation with the semivariogram
- Estimate the semivariogram from point observations
- Use the semivariogram for spatial interpolation with ordinary kriging
- Extend ordinary kriging to regression kriging

ISRIC **World Soil Information**

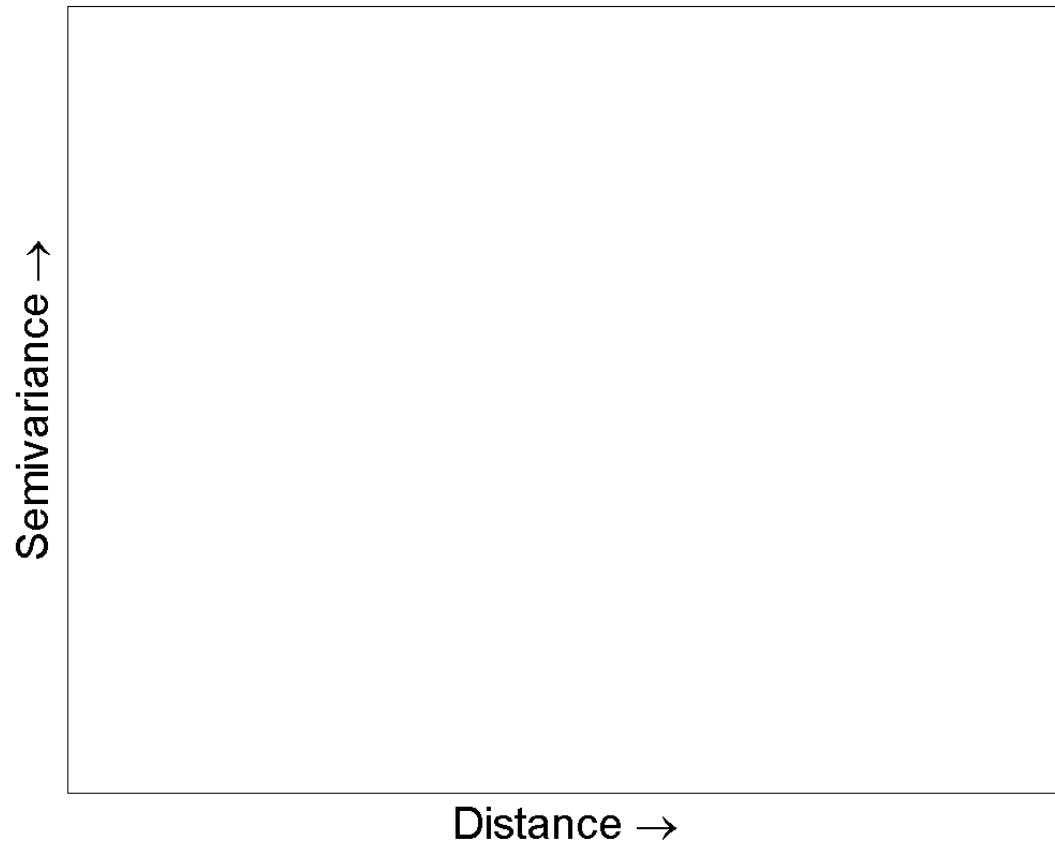# What we all know: soils and soil properties vary in space

Spatial variation can be quantified using the so-called **semivariance**, this is half the expected squared difference between the values of the variable of interest at two locations

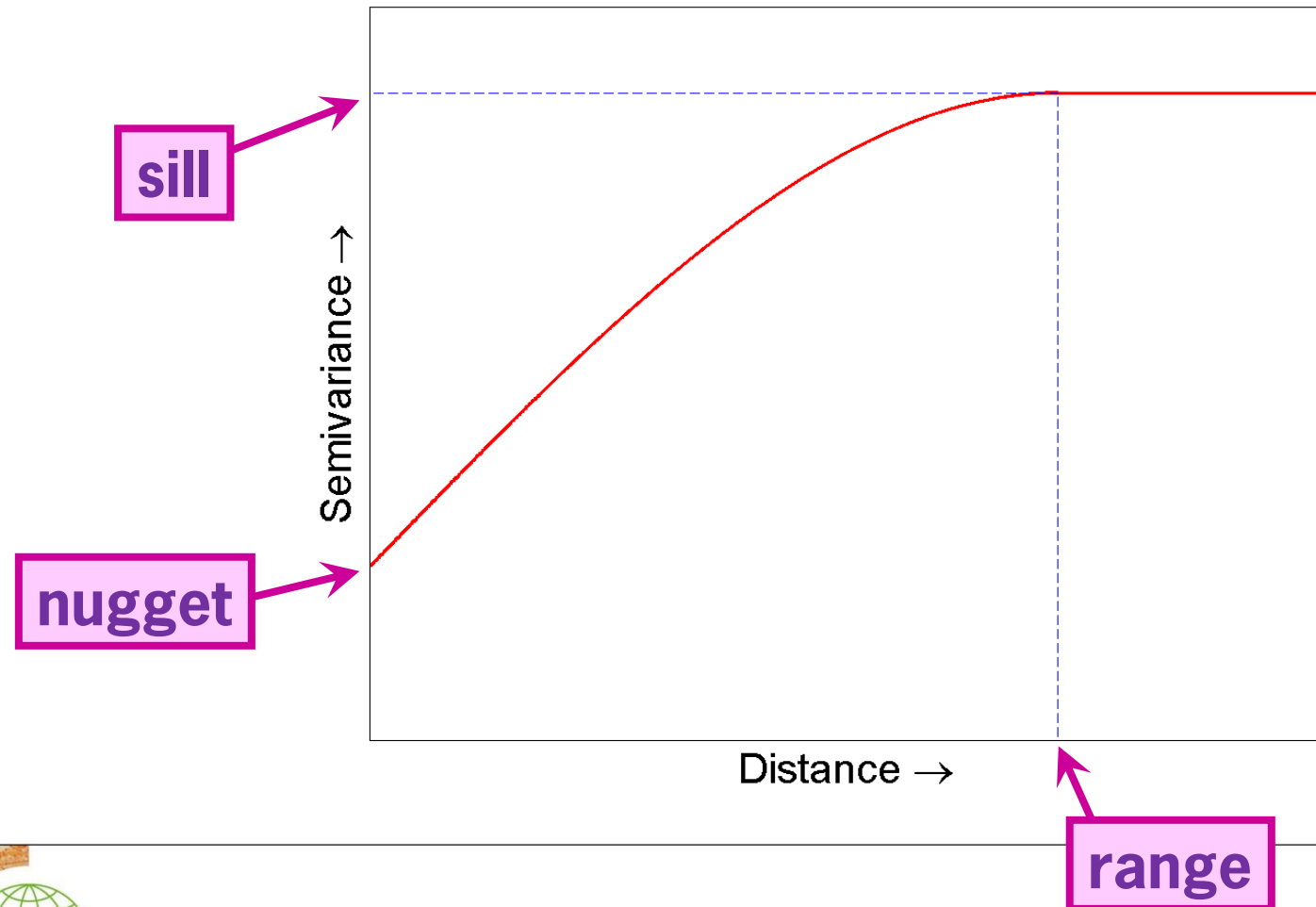$$\gamma(h) = \frac{1}{2} E\left[\left(Z(x) - Z(x+h)\right)^2\right]$$

*measurement at location* x

*measurement at location* x+h

# Plot of semivariance as a function of the distance is called a **(semi)variogram**
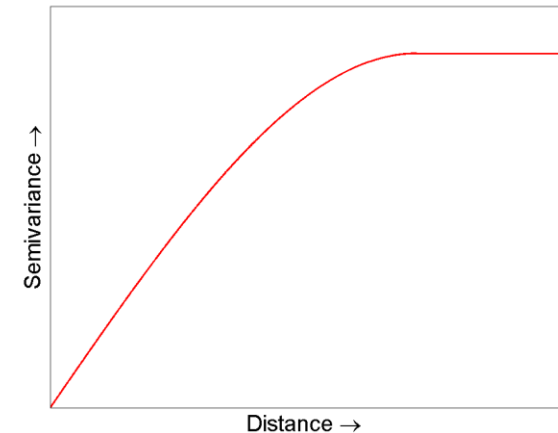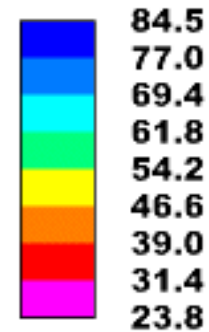
# Typical shape of the semivariogram, with parameters
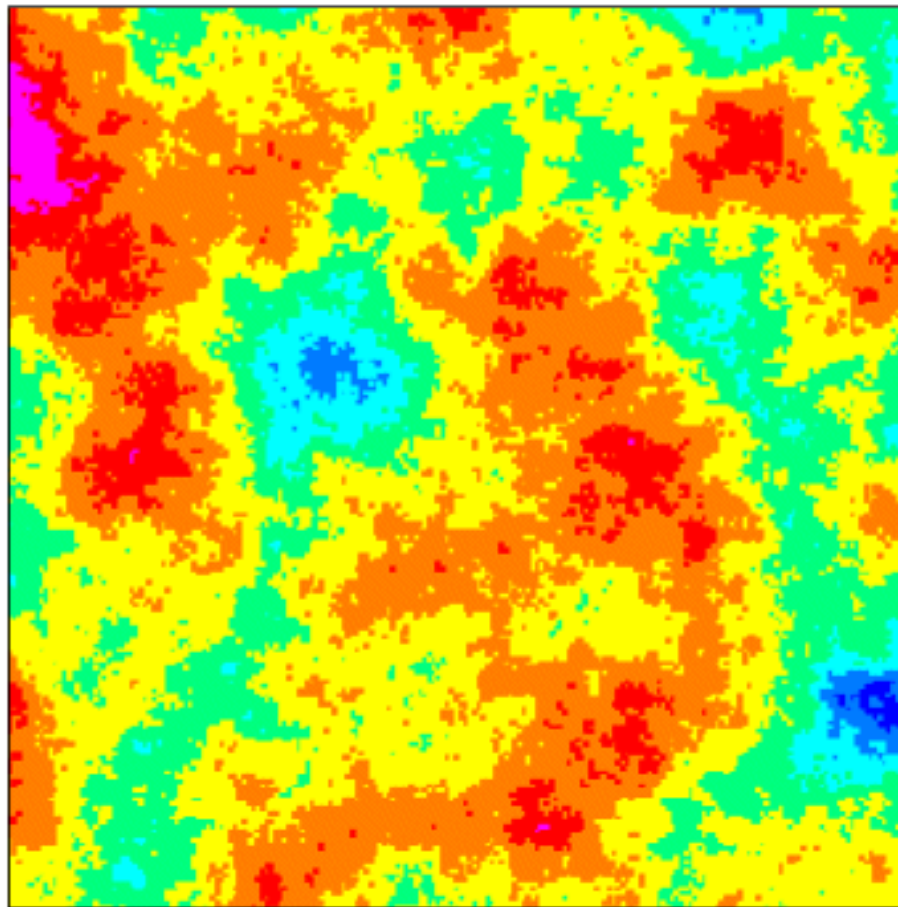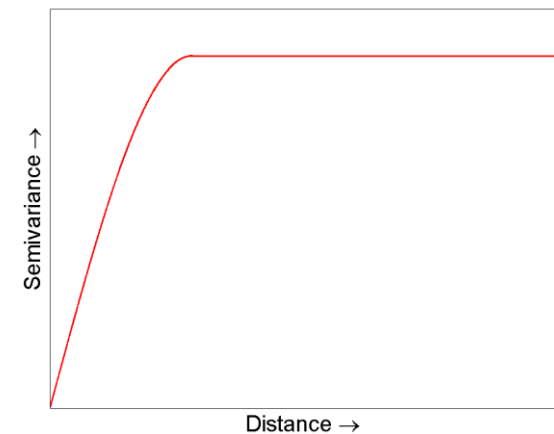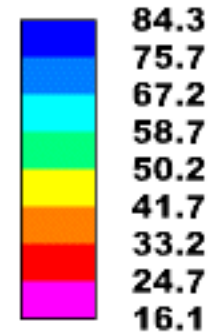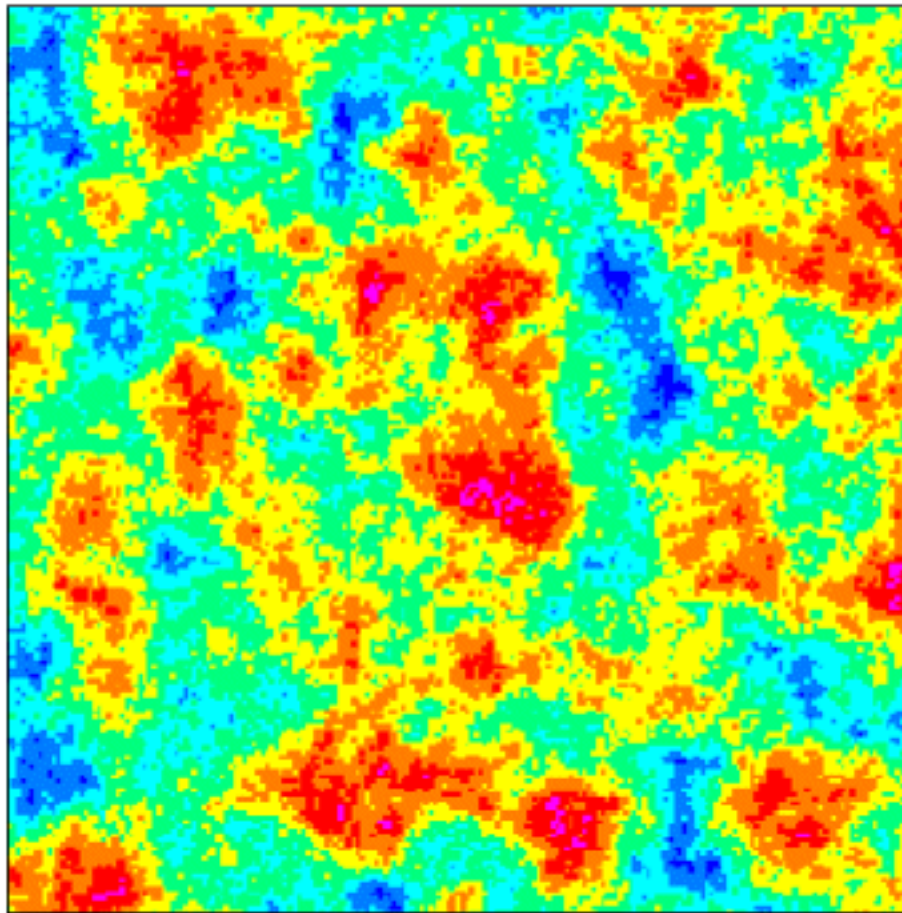
# Interpretation of semivariogram parameters

- **Nugget** = measurement errors and short-distance spatial variation
- **Sill** = variance of the variable of interest
- **Range** = distance up to which there is spatial correlation



ISRIC **World Soil Information**

# Example 1: a possible 'reality'

# Example 2: another reality

# Example 3: third reality

# Example 4: fourth reality

# When presenting the semivariogram, two assumptions were implicitly made:

1.  The semivariance of $Z(x)$ and $Z(x+h)$ only depends on the distance $h$ and not on the locations $x$ and $x+h$ (**stationarity** assumption)
2.  The semivariance is a function of the length of $h$, not of its direction (**isotropy** assumption)

These assumptions are not always realistic and can be relaxed, but today we won't go into that

ISRIC  World Soil Information

# Estimation of the semivariogram from point observations, take pH data Europe as an example



pH observations A horizon

World Soil Information

# Structural analysis: estimate semivariogram from observations

- Suppose there are $n$ observations (in this example $n=2582$)
- This yields $\frac{1}{2} \times n \times (n-1)$ pairs of observations
- Each pair of observations $\{z(x_i), z(x_j)\}$ provides information about the semivariance over distance $|x_i - x_j|$ (by computing $\frac{1}{2} \times (z(x_i) - z(x_j))^2$)
- Presented in a graph: semivariogram cloud

ISRIC World Soil Information

# Semivariogram cloud pH data



- How many blue crosses are in this figure?

# Averaging over 'lags' (intervals) gives experimental semivariogram

# Last step of structural analysis: fit a function through the experimental semivariogram

1. Choose a function shape, common choices:

spherical :

$$\gamma(h) = 0 \quad \text{for } h = 0$$

$$\gamma(h) = c_0 + c \cdot \{\frac{3}{2} \cdot \frac{h}{a} - \frac{1}{2} \cdot \left(\frac{h}{a}\right)^3\} \quad \text{for } 0 < h \leq a$$

$$\gamma(h) = c_0 + c \quad \text{for } h > a$$

linear :

$$\gamma(h) = 0 \quad \text{for } h = 0$$

$$\gamma(h) = c_0 + b \cdot h \quad \text{for } h > 0$$

exponential :

$$\gamma(h) = 0 \quad \text{for } h = 0$$

$$\gamma(h) = c_0 + c \cdot (1 - e^{-\frac{h}{a}}) \quad \text{for } h > 0$$

Gaussian :

$$\gamma(h) = 0 \quad \text{for } h = 0$$

$$\gamma(h) = c_0 + c \cdot (1 - e^{-\left(\frac{h}{a}\right)^2}) \quad \text{for } h > 0$$

2. Estimate parameters of the chosen shape (e.g. using weighted least squares fitting)

ISRIC World Soil Information

# Resulting variogram model for pH example:



Shape = Spherical

Nugget = 0.57

Sill = 1.81

Range = 812 km

- Does it agree with our hypothesis?

World Soil Information

# Geostatistical interpolation: Kriging

- Introduced in the 1950s by Daniel Krige: mining engineer from South-Africa
- Kriging comes in many forms, we focus on Ordinary Kriging but will also look at Regression Kriging
- Principle: prediction at a location is a linear combination of observations nearby
- The weight that is given to each observation depends on the degree of (spatial) correlation: the semivariogram plays an important role

ISRIC World Soil Information

# Ordinary Kriging

Predict $Z(x_0)$ at unobserved location $s_0$ using observations $Z(x_i)$, $i=1,\ldots,n$ as follows:

$$\hat{Z}_{OK}(x_0) = \sum_{i=1}^{n} \lambda_i \cdot Z(x_i)$$

**Kriging weight**

**ISRIC** **World Soil Information**

# How to choose the kriging weights?

- What criterion would you recommend?

World Soil Information

# How to choose the kriging weights?

- What criterion would you recommend?



- What would be nice properties of the kriging prediction error $\hat{Z}(x_0) - Z(x_0)$?

# Which probability distribution of $\hat{Z}(x_0) - Z(x_0)$ do you prefer?

# Computation of kriging weights

Minimise the expected squared prediction error, in other words choose the coefficients $\lambda_i$ such that:

$$E[(\hat{Z}(x_0) - Z(x_0))^2]$$

is as small as possible, under the unbiasedness condition:

$$\sum_{1=1}^{n} \lambda_i = 1$$

ISRIC **World Soil Information**

# Solution in case of Ordinary Kriging:

**Lagrange parameter**

$$\sum_{j=1}^{n} \lambda_j \cdot \gamma(|x_i - x_j|) + \varphi = \gamma(|x_i - x_0|) \quad for \ all \ i = 1, ..., n$$

$$In \ addition \quad \sum_{i=1}^{n} \lambda_i = 1$$

**ISRIC** **World Soil Information**

Which observation gets the largest weight? Which the smallest? What would we get in case of inverse distance interpolation?



$x_1$

$x_2$

$x_0$

$x_3$

$x_4$

ISRIC **World Soil Information**

# Application to European pH data



pH observations A horizon

Legend:
- [2.9,4.3]
- (4.3,5.7]
- (5.7,7.1]
- (7.1,8.5]
- (8.5,9.9]

Semivariogam pH A horizon

# OK prediction yields a smoothed representation of reality



pH observations A horizon

[2.9,4.3]
(4.3,5.7]
(5.7,7.1]
(7.1,8.5]
(8.5,9.9]

OK prediction pH A horizon

**ISRIC** World Soil Information

# The kriging interpolation error is quantified by the Ordinary Kriging variance:

semivariance at distance $x_i - x_0$

$$\sigma^2_{OK}(x_0) = E\left[(Z(x_0) - \hat{Z}(x_0))^2\right] = \sum_{i=1}^{n} \lambda_i \cdot \gamma(x_i - x_0) + \varphi$$

This shows that the interpolation error is small when there is strong spatial correlation and/or when there are many observations in the local neighbourhood

ISRIC World Soil Information

# Map of OK standard deviations ($=\sqrt{{\sigma_{OK}}^2}$) shows observation density



pH observations A horizon

OK standard deviation pH A horizon

# Regression kriging

$$Z(x) = m(x) + \varepsilon(x)$$

dependent, target variable

trend, explanatory part

stochastic residual, unexplanatory part, can be spatially correlated!

Unlike ordinary kriging, in regression kriging the trend is no longer constant but a function of 'explanatory' variables, for example:

$$soil\ depth(x) = \begin{aligned} &\beta_0 + \beta_1 \cdot elevation(x) + \beta_2 \cdot slope\ angle(x) \\ &+ \beta_3 \cdot vegetation\ density(x) \\ &+ \beta_4 \cdot upstream\ area(x) \end{aligned} + residual(x)$$

# Regression kriging algorithm

1. select explanatory variables and fit regression model (estimate regression coefficients)

2. compute residuals (by subtracting the fitted trend from the observations) at observation locations and compute from them a semivariogram

3. apply the regression model to all unobserved locations (usually a grid)

4. krige the residuals

5. add up the results of steps 3 and 4

ISRIC **World Soil Information**

# Example from Hengl et al. (Geoderma, 2004): predicting soil depth for a 50 × 50 km area in Croatia

# Results using four interpolation methods



observations

soil map only predictor

regression only

ordinary kriging

regression kriging

**ISRIC** World Soil Information

# Validation on 35 independent observations

| | Mean error [cm] | Root mean squared error [cm] |
|---|---|---|
| Soil map | 1.42 | 9.1 |
| Ordinary kriging | 0.69 | 8.5 |
| Multiple regression | 1.69 | 8.8 |
| Regression kriging | 0.15 | 6.8 |

ISRIC **World Soil Information**

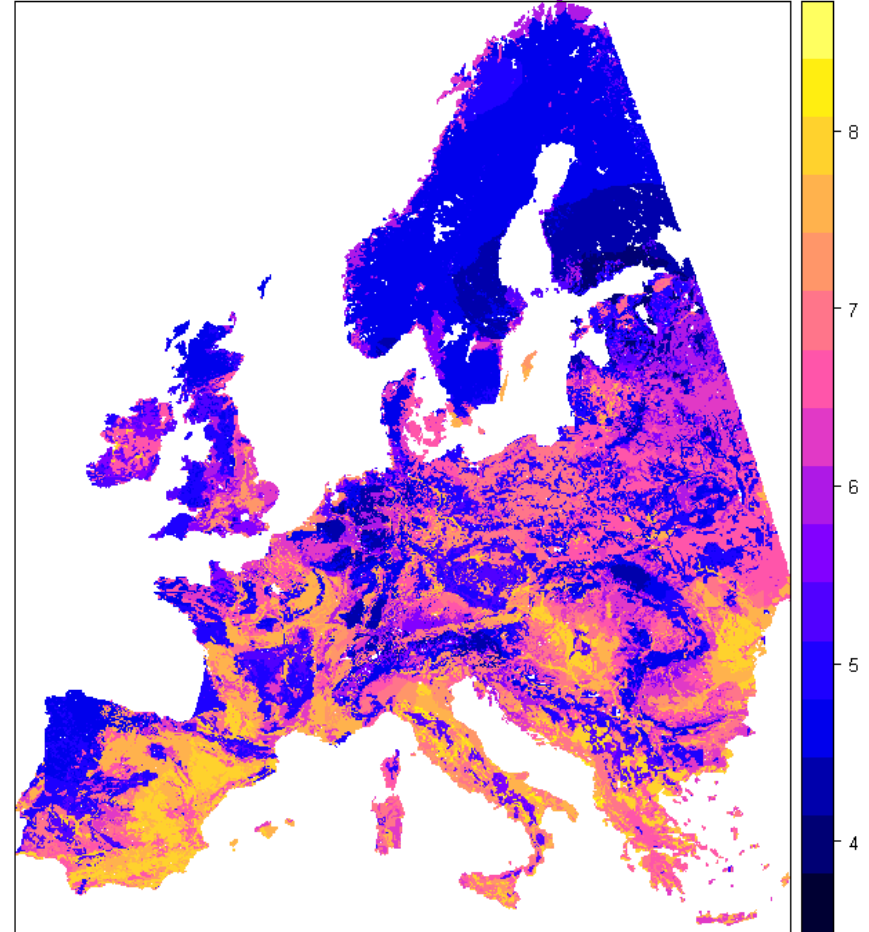# Regression kriging also yields more realistic results for pH in Europe

## Ordinary Kriging

## Regression Kriging

# Which sampling design is best?

- Variogram estimation has different requirements than kriging: quantification of <span style="color:red">short-distance spatial variation</span> versus <span style="color:red">uniform spread</span>, pragmatic solution:



o gridpoints

■ short distance point

# If you really want the optimum design: given the semivariogram, it can be obtained with numerical techniques such as spatial simulated annealing



Fig. 12. An a priori optimised sampling scheme for anisotropic sand percentage.

0.00  0.20  0.40  0.60  0.80  km

**ISRIC** World Soil Information

# Optimization in feature space

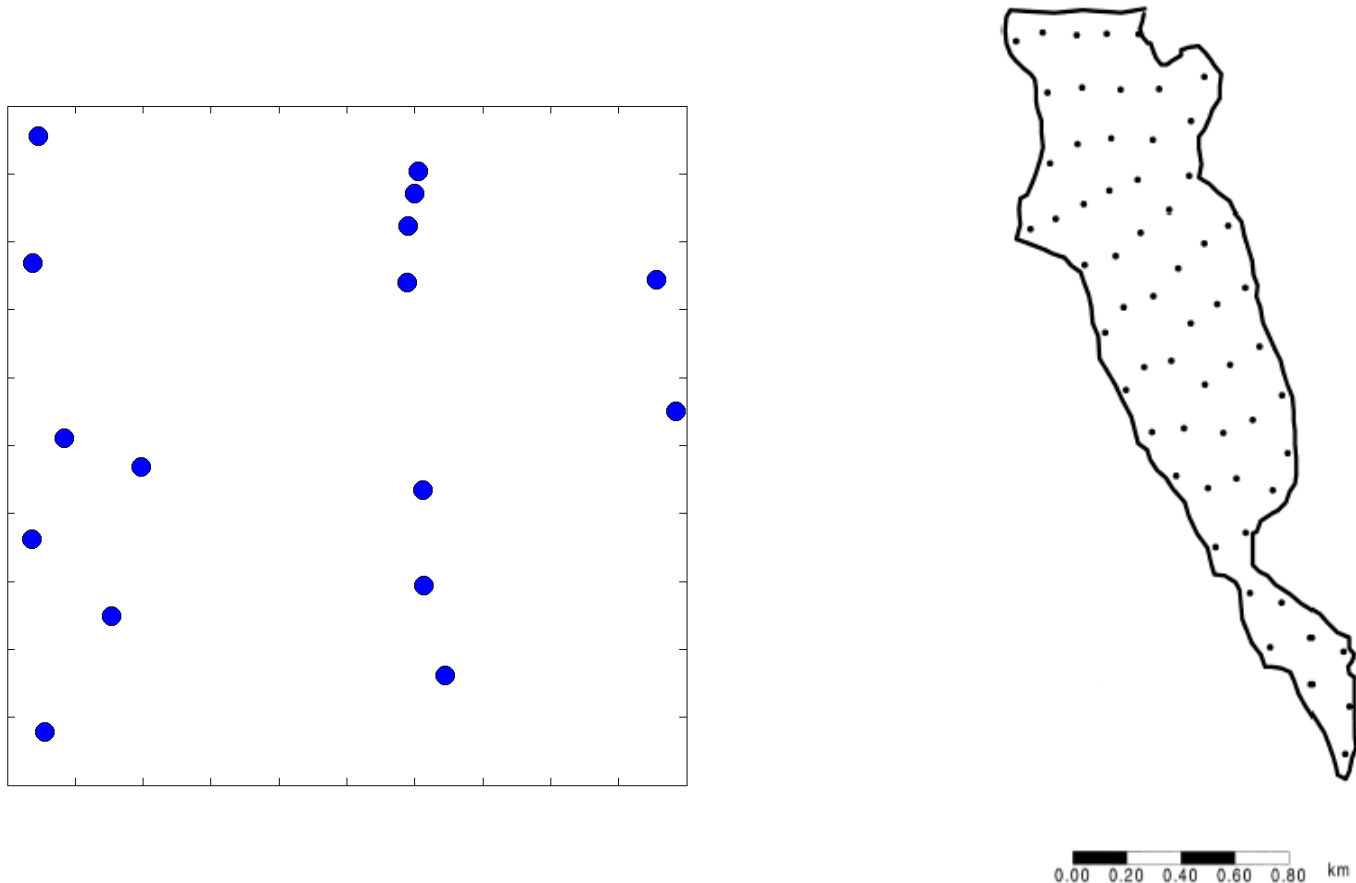$\sigma_{RK}{}^2(x) =$

spatial

dependent →

expla

# A conditioned Latin hypercube method for sampling in the presence of ancillary information

Budiman Minasny[*], Alex B. McBratney

*Australian Centre for Precision Agriculture, Faculty of Agriculture, Food and Natural Resources, The University of Sydney, Australia*

## Abstract

This paper presents the conditioned Latin hypercube as a sampling strategy of an area with prior information represented as exhaustive ancillary data. Latin hypercube sampling (LHS) is a stratified random procedure that provides an efficient way of sampling variables from their multivariate distributions. It provides a full coverage of the range of each variable by maximally stratifying the marginal distribution. For conditioned Latin hypercube sampling (cLHS) the problem is: given $N$ sites with ancillary variables $(X)$, select $x$ a sub-sample of size $n$ $(n \ll N)$ in order that $x$ forms a Latin hypercube, or the multivariate distribution of $X$ is maximally stratified. This paper presents the cLHS method with a search algorithm based on heuristic rules combined with an annealing schedule. The method is illustrated with a simple 3-D example and an application in digital soil mapping of part of the Hunter Valley of New South Wales, Australia. Comparison is made with other methods: random sampling, and equal spatial strata. The results show that the cLHS is the most effective way to replicate the distribution of the variables.

ISRIC