

Applied geostatistics

Lecture 5 – Spatial prediction from point samples (Part 2)

D G Rossiter

University of Twente.

Faculty of Geo-information Science & Earth Observation (ITC) – August 2014

Section of Soil & Crop Sciences, Cornell University September 2014–

February 6, 2016

Copyright © 2016 D G Rossiter

All rights reserved. Reproduction and dissemination of this work as a whole (not parts) freely permitted if this original copyright notice is included. Sale or placement on a web site where payment must be made to access this document is strictly prohibited. To adapt or translate please contact the author (dgr2@cornell.edu).



Topics for this lecture

1. Derivation of the Ordinary Kriging (OK) system: (1) *regression* (2) minimization
2. *Simple Kriging* (SK)
3. Block Kriging
4. Universal Kriging (UK)
5. Derivation of the Universal Kriging system: (1) *regression* (2) minimization
6. *Kriging transformed variables*
7. *Kriging with External Drift* (KED) and *Regression Kriging* (RK)
8. *Stratified Kriging* (StK)
9. *Co-Kriging* (Co-K)

The topics written *in italic script* are optional at first reading.

Commentary

This lecture has two purposes:

1. Show how kriging is in some sense an “optimal” predictor;
2. Present various kriging variants

These many variants all are applicable in certain situations, which are explained with each variant.

Deriving the kriging system

In the previous lecture we saw how to apply a **model of local spatial dependence** (i.e. a variogram model) to **prediction by kriging**.

To avoid information overload, we deferred discussing the **kriging equations**, and in particular in what sense kriging is an **optimal** local predictor.

Note: It is not necessary to understand this topic completely in order to correctly apply kriging. The derivation is necessarily mathematical and in places requires knowledge of matrix algebra or differential calculus. Still, everyone who uses kriging should be exposed to this at least once.

Two approaches

There are two approaches to this derivation:

Regression As a special case of **weighted least-squares** prediction in the **generalized linear model** with orthogonal projections in linear algebra

Minimization **Minimizing the kriging prediction variance** with calculus

Approach (1) is mathematically more elegant and provides a clear link to well-established linear modelling theory, so we present it as the main derivation.

Approach (2) is an application of standard minimization methods from differential calculus; but is not so transparent, because of the use of LaGrange multipliers.

Topic: Regression derivation of the Ordinary Kriging equations

Here we show how to derive and solve the kriging equations in a uniform framework for linear modelling (“regression”) and kriging.

This approach has four steps:

1. Model the **spatial structure**, e.g. the covariance function or semivariogram function;
2. Estimate the **spatial mean** (not necessary in Simple Kriging (SK));
3. Set up a **kriging system** to **minimize the prediction variance**;
4. Compute the **kriging weights**

Step (1) has been discussed in a previous topic.

Note: If you are more comfortable with differential calculus than with linear algebra, the minimization derivation (next section) will likely be more accessible.

Commentary

This topic requires a knowledge of **matrix algebra** and some familiarity with **general linear models** (e.g. weighted least-squares).

Step 2: estimate the spatial mean

The **spatial mean** of the variable z to be predicted, over the study area A , with area $|A|$ is:

$$\hat{\mu}(z) = \frac{1}{|A|} \int_A z$$

In practice this is **discretized** by summing over some fine **grid**:

$$\hat{\mu}(z) = \frac{1}{x_{\max} y_{\max}} \sum_{i=1 \dots x_{\max}, j=1 \dots y_{\max}} z_{i,j}$$

where $z_{i,j}$ is the value of the variable at grid location (i, j) .

But in general we do not have measurements at all locations! That's what we want to find out. So this equation can't be applied. We need to estimate the spatial mean from **sparse observations** (i.e. our sample points)

The spatial mean is not the average!

Problem: because of **spatial autocorrelation**, it is not correct simply to **average** the observations to obtain the mean.

If there is any spatial dependence, in general:

Spatial mean \neq **Average** of the observations

To check your understanding . . .

Q1 : *Consider a regular grid over some area; measure the data values and compute the ordinary average.*

If there is spatial dependence, what would happen to this average if we now make many observations near one of the grid points (i.e. a cluster), and not the others?

Jump to A1 •

The spatial mean

This is solved by weighted least squares, taking into account spatial correlation:

$$\hat{\mu} = (\mathbf{1}^T \cdot \mathbf{C}^{-1} \cdot \mathbf{1})^{-1} \cdot (\mathbf{1}^T \cdot \mathbf{C}^{-1} \cdot \mathbf{z})$$

- $\mathbf{1}$: a $n \times 1$ column vector of n 1's; so $\mathbf{1}^T$: a $1 \times n$ row vector of n 1's
 - * In Universal Kriging (see below) this will be generalized to a design matrix \mathbf{Q} ; here just a vector of 1's to estimate the mean
- \mathbf{C} : the covariance matrix ($n \times n$) among known points;
 - * Note that $C(\mathbf{0}) \equiv 1$ so all diagonals are 1.
- \mathbf{z} : the $n \times 1$ column vector of the n known data values

This is a special case of Generalized Least Squares (GLS).

To check your understanding ...

Q2 : *How is the covariance matrix \mathbf{C} computed?*

Jump to A2 •

Breaking this down ...

Two **quadratic forms**; both are $1 \times n \cdot n \times n \cdot n \times 1$ and so end up as scalars; multiply two scalars to get the spatial mean:

$$\begin{aligned} & (\mathbf{1}^T \cdot \mathbf{C}^{-1} \cdot \mathbf{1})^{-1} \\ & (\mathbf{1}^T \cdot \mathbf{C}^{-1} \cdot \mathbf{z}) \end{aligned}$$

Note: Let quadratic form $(\mathbf{1}^T \cdot \mathbf{C}^{-1} \cdot \mathbf{1})^{-1} = \mathbf{V}$ in subsequent formulas.

If there is no spatial correlation, \mathbf{C} reduces to \mathbf{I} : all diagonals are 1, all off-diagonals are 0, the inverse $\mathbf{I}^{-1} = \mathbf{I}$, and this reduces to Ordinary Least Squares (OLS) estimation.

In this case (design matrix $\mathbf{Q} = \mathbf{1}$), the OLS estimation of the spatial mean becomes the **arithmetic average**: the first quadratic term is $1/(\sum 1) = 1/n$ and the second $\sum z$; their product is $1/n \cdot \sum z$

Spatial mean computed with semivariances

Recall: the **semivariance** is the deviation of the covariance at some separation \mathbf{h} from the total variance:

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

i.e.

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

But $C(\mathbf{0})$ is constant (1) in the covariance functions; further, both quadratic forms include the \mathbf{C} matrix, so using its negative (plus a constant term), e.g. Γ , does not change the solution.

In fact, we can replace C by its difference from an arbitrary scalar α (element-wise):

$$\mathbf{C} \rightarrow \alpha \cdot \mathbf{1} \cdot \mathbf{1}^T - \mathbf{C}$$

So **semivariances** may be used rather than **covariances** in the formula for the spatial mean.

Step 3: Kriging prediction

Once we know the spatial mean, a kriging system can be set up without any additional constraints.

The special property of this system is that it is **BLUP**: “**B**est **L**inear **U**nbiased **P**redictor”, given the modelled covariance structure.

Formula for kriging prediction variance

By definition, this is:

$$\text{Var}(\hat{Z}_0 - Z_0)$$

where \hat{Z}_0 is the predicted value and Z_0 is the (unknown) true value.

For *any* **weighted average**, with weights λ_i , this is:

$$\text{Var}(\sum \lambda_i \cdot Z_i - Z_0)$$

Kriging selects the weights λ to minimize this expression.

Continued . . .

To make the division of weights clear, define two vectors of length $n + 1$:

$$v = \begin{pmatrix} \lambda \\ -1 \end{pmatrix} \quad Z = \begin{pmatrix} Z \\ Z_0 \end{pmatrix}$$

where the first n elements relate to the observation points and the last element to the prediction point. Then

$$\begin{aligned} \text{Var}(\sum \lambda_i \cdot Z_i - Z_0) &= v^T \text{var}(Z) v \\ &= v^T \begin{pmatrix} C & c_0 \\ c_0^T & c_{00} \end{pmatrix} v \\ &= \lambda^T C \lambda - 2\lambda^T c_0 + c_{00} \end{aligned}$$

The key point here is that the variance-covariance matrix $\text{var}(Z)$ is broken down into submatrices: C for the covariances between sample points, c_0 for the covariance of each sample point with the prediction point, and c_{00} for the variance at a point (the nugget).

Formula for kriging prediction

$$\hat{Z}_{OK} = \hat{\mu} + \mathbf{c}_0^T \cdot \mathbf{C}^{-1} \cdot (\mathbf{z} - \hat{\mu}\mathbf{1})$$

- \hat{Z}_{OK} is the **prediction** at the prediction point
- $\hat{\mu}$ is the **spatial mean** computed in the previous step
- \mathbf{c}_0 is a $n \times 1$ column vector of the **covariance** between **each sample point** and the point to be **predicted**
- \mathbf{C} : the **covariance matrix** ($n \times n$) among known points
- \mathbf{z} : the $n \times 1$ column vector of the n **known data values**
- $\mathbf{1}$: a $n \times 1$ column vector of n 1's, so that $\hat{\mu}\mathbf{1}$ is a column vector of the means, and $(\mathbf{z} - \hat{\mu}\mathbf{1})$ is a column vector of the **residuals** from the spatial mean

To check your understanding ...

Q3 : *How is the covariance vector \mathbf{c}_0 computed?*

Jump to A3 •

Breaking this down ...

- Note that the spatial mean $\hat{\mu}$ is subtracted from each data value; so the right summand

$$\mathbf{c}_0^T \cdot \mathbf{C}^{-1} \cdot (\mathbf{z} - \hat{\mu}\mathbf{1})$$

is the **deviation** from the spatial mean at the prediction point.

- The spatial mean is added back in the left summand.
- If there is no spatial dependence, $\mathbf{c}_0 = \mathbf{0}$ and the prediction is just the spatial mean (which would then just be the average).

Step 4: Kriging weights

The OK prediction equation just derived does not explicitly give the weight to each observation; the kriging system can be used directly without computing these.

But we would often like to know the weights, to see the relative importance of each observation; also it may be more efficient to compute these and then the prediction as the weighted average.

Recall: the vector of kriging **weights**, is what to multiply each observation by in the weighted sum (prediction).

We get this by collecting all the terms that multiply the observed values \mathbf{z} in the OK prediction equation:

$$\lambda^T = \mathbf{c}_0^T \cdot \mathbf{C}^{-1} - \mathbf{1}^T \cdot \mathbf{V} \mathbf{1}^T \cdot \mathbf{C}^{-1}$$

Kriging variance without weights

Above the prediction variance was computed for *any* set of weights; now we've selected weights to minimize the variance. Then the variance can be expressed without explicitly showing the weights.

We do this by substituting the expression for weights (previous step) into the expression for prediction variance (Step 3), to obtain:

$$\text{Var}(\hat{Z}_0 - Z_0) = c_{00} - \mathbf{c}_0^T \cdot \mathbf{C}^{-1} \cdot \mathbf{c}_0 + \mathbf{x}_a^T \cdot \mathbf{V} \cdot \mathbf{x}_a$$

where:

- c_{00} is the nugget covariance (at separation 0).
- $\mathbf{x}_a = \mathbf{1} - (\mathbf{c}_0^T \cdot \mathbf{C}^{-1} \cdot \mathbf{1})$

Note: For **block kriging**, replace c_{00} with c_{BB} , the average within-block variance (which will in general be smaller than the at-point variance); replace \mathbf{c}_0 and \mathbf{C} with block-to-block covariances (see topic “Block Kriging”).

Kriging variance in terms of semivariances

Because of the relation

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

the kriging variance can also be expressed as:

$$\text{Var}(\hat{Z}_0 - Z_0) = c_{00} - \gamma_0^T \cdot \Gamma^{-1} \cdot \gamma_0 - x_a/\mathbf{V}$$

where:

- Γ is the matrix of semivariances between sample points
- γ_0 is the vector of semivariances between sample points and the point to be predicted

Note the change of sign, and that c_{00} is now implicit.

Topic: Minimization derivation of the Ordinary Kriging equations

This *optional* topic presents another approach to deriving the kriging equations.

If you are more comfortable with differential calculus than with linear algebra, this derivation will likely be more accessible.

Commentary

In the detailed derivation that follows, keep your eye on the **big picture**:

- We want an **optimal** linear predictor, i.e. we want to compute **the best possible** weights for each sample point, to make each prediction;
- “Optimal” depends on some **objective function** which can be **minimized** with the best weights;
- We choose the **variance of the prediction** as the **objective function**; i.e. we want to **minimize** the **uncertainty** of the prediction.
- Given all this, we can derive a set of linear equations to derive the weights; this is known as the **kriging system**.
- The kriging system takes into account the **relative spatial positions** of each sample point, and their positions with respect to the prediction point.
- It does this with a **model of spatial covariance**.

Prediction Variance

For **any** predictor (not just the kriging predictor):

- The **prediction** $\hat{Z}(\mathbf{x}_0)$ at a given location \mathbf{x}_0 may be compared to the **true** value $Z(\mathbf{x}_0)$; note the “hat” symbol to indicate an **estimated** value rather than a **measured** one.
- Even though we don’t know the true (measured) value, we can write the expression for the **kriging variance**
- This is defined as the **expected value** of the **squared difference** between the **estimate** and the (unknown) **true value**:

$$\sigma^2(Z(\mathbf{x}_0)) \equiv E[\{\hat{Z}(\mathbf{x}_0) - Z(\mathbf{x}_0)\}^2]$$

- If we can express this in some **computable form** (i.e. without the unknown true value) we can use it as an **optimality criterion**

Derivation of the kriging variance for OK

(Derivation based on P K Kitanidis, *Introduction to geostatistics: applications to hydrogeology*, Cambridge University Press, 1997; §3.9)

1. In OK, the estimated value is a **linear combination** of data values \mathbf{x}_i , with weights λ_i derived from the **kriging system**:

$$\hat{z}(\mathbf{x}_0) = \sum_{i=1}^N \lambda_i z(\mathbf{x}_i)$$

2. We don't yet know the weights λ but we do know how to compute the **estimated value** given the weights.
3. So, we can re-write the kriging variance with this weighted sum:

$$\sigma^2(Z(\mathbf{x}_0)) = E\left[\left\{\sum_{i=1}^N \lambda_i z(\mathbf{x}_i) - Z(\mathbf{x}_0)\right\}^2\right]$$

Expanding into parts

4. Add and subtract the (unknown but stationary) mean μ :

$$\sigma^2(z(\mathbf{x}_0)) = E\left[\left\{\sum_{i=1}^N \lambda_i(z(\mathbf{x}_i) - \mu) - (Z(\mathbf{x}_0) - \mu)\right\}^2\right]$$

5. Expand the square:

$$\sigma^2(Z(\mathbf{x}_0)) = E\left[\left(\sum_{i=1}^N \lambda_i z(\mathbf{x}_i) - \mu\right)^2 - 2 \sum_{i=1}^N \lambda_i (z(\mathbf{x}_i) - \mu)(Z(\mathbf{x}_0) - \mu) + (Z(\mathbf{x}_0) - \mu)^2\right]$$

6. Replace the squared single summation (first term) by a **double summation**, i.e. with separate indices for the two parts of the square:

$$\left(\sum_{i=1}^N \lambda_i z(\mathbf{x}_i) - \mu\right)^2 = \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j (z(\mathbf{x}_i) - \mu)(z(\mathbf{x}_j) - \mu)$$

Bring expectations into each term

7. Bring the expectation into each term (n.b. expectation of a sum is the sum of expectations):

$$\begin{aligned}\sigma^2(Z(\mathbf{x}_0)) &= \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j E[(z(\mathbf{x}_i) - \mu)(z(\mathbf{x}_j) - \mu)] \\ &\quad - 2 \sum_{i=1}^N \lambda_i E[(z(\mathbf{x}_i) - \mu)(Z(\mathbf{x}_0) - \mu)] + E[(Z(\mathbf{x}_0) - \mu)^2]\end{aligned}$$

From expectations to covariances

8. Now, the three **expectations** in the previous expression are the definitions of **covariance** or **variance**:

- (a) $E[(z(\mathbf{x}_i) - \mu)(z(\mathbf{x}_j) - \mu)]$: **covariance** between **two sample points**
- (b) $E[(z(\mathbf{x}_i) - \mu)(Z(\mathbf{x}_0) - \mu)]$: **covariance** between **one sample point** and the **prediction point**
- (c) $E[(Z(\mathbf{x}_0) - \mu)^2]$: **variance** at the **prediction point**

9. So, replace the **expectations** with **covariances** and **variances**:

$$\begin{aligned}\sigma^2(Z(\mathbf{x}_0)) &= \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j \text{Cov}(z(\mathbf{x}_i), z(\mathbf{x}_j)) \\ &\quad - 2 \sum_{i=1}^N \lambda_i \text{Cov}(z(\mathbf{x}_i), Z(\mathbf{x}_0)) + \text{Var}(Z(\mathbf{x}_0))\end{aligned}$$

How can we evaluate this expression?

- **Problem 1:** how do we know the **covariances** between any two points?
 - * Answer: by applying a **covariance function** which **only depends on spatial separation** between them.
 - Authorized functions e.g. spherical, exponential, ...
- **Problem 2:** how do we find the correct covariance function?
 - * Answer: this is the subject of **variogram analysis**
- **Problem 3:** how do we know the **variance** at any point?
 - * Answer: The actual value doesn't matter (it will be eliminated in the following algebra) but it **must be the same at all points**: this is the assumption of **second-order stationarity**.

Stationarity

- This is a term for **restrictions on the nature of spatial variation** that are required for OK to be correct
- **First-order** stationarity: the **expected values** (mean) at all locations in the field are the **same**:

$$E[Z(\mathbf{x}_i)] = \mu, \forall \mathbf{x}_i \in R$$

- **Second-order** stationarity:
 1. The **variance** at any point is **finite** and the **same at all locations** in the field
 2. The **covariance structure** depends only on **separation** between point pairs
- Note: it is easy to overcome problems with first-order stationarity (changing expected value across the field) but not second-order (changing variance and covariance)

Commentary

The concept of stationarity is often confusing, because stationarity refers to **expected** values, variances, or co-variances, rather than **observed** values.

In particular, of course the actual values change over the field! That is exactly what we want to use to predict at unsampled points. First-order stationarity just says that **before we sampled**, the **expected** value at all locations was the same.

That is, we **assume** the values result from a **spatially-correlated process** with a **constant mean** – **not** constant values!

Once we have some sample values, these influence the probability of finding values at other points, because of **spatial covariance**.

Unbiasedness

An **unbiased** estimate is one where the **expectation** of the **estimate** equals the expectation of the **true** (unknown) value: $E[\hat{Z}(\mathbf{x}_0)] \equiv E[Z(\mathbf{x}_0)]$

For OK, we have assumed $E[Z(\mathbf{x}_0)]$ is some unknown but **constant** μ . We have also decided to estimate $E[\hat{Z}(\mathbf{x}_0)]$ as a weighted sum, i.e. linear combination; see Step 1 previous slide. So we must have:

$$E[\hat{Z}(\mathbf{x}_0)] = \sum_{i=1}^N \lambda_i E[z(\mathbf{x}_i)] = \sum_{i=1}^N \lambda_i \mu = \mu \sum_{i=1}^N \lambda_i$$

because all the expected values are the same μ by **first-order stationarity**.

Since $E[\hat{Z}(\mathbf{x}_0)] = \mu$ (unbiasedness), we must have $\sum_{i=1}^N \lambda_i = 1$.

This will be an additional restriction in the kriging system.

From point-pairs to separation vectors

As written above, the expressions are huge! The covariances between all point-pairs must be determined separately.

However, because of **second-order stationarity**, we assume that the covariances between any two points depend **only on their separation** and a single **covariance function**.

So rather than try to compute all the covariances, we just need to know this function, then we can apply it to any point-pair, just by knowing their separation.

We continue the derivation ...

Replace point-pairs with separation vectors

10. Substitute the **covariance function** of **separation \mathbf{h}** into the expression; note that at separation zero (a point) this is the variance:

$$\sigma^2(Z(\mathbf{x}_0)) = \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j \text{Cov}(\mathbf{h}(i, j)) - 2 \sum_{i=1}^N \lambda_i \text{Cov}(\mathbf{h}(i, 0)) + \text{Cov}(\mathbf{0})$$

- $\mathbf{h}(i, 0)$ is the separation between sample point \mathbf{x}_i and the point to be predicted \mathbf{x}_0 .
- $\mathbf{h}(i, j)$ is the separation between two sample points \mathbf{x}_i and \mathbf{x}_j .

From covariances to semi-variances

11. Replace covariances by **semivariances**, using the relation $\text{Cov}(\mathbf{h}) = \text{Cov}(\mathbf{0}) - \gamma(\mathbf{h})$:

$$\sigma^2(Z(\mathbf{x}_0)) = - \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j \gamma(\mathbf{h}(i, j)) + 2 \sum_{i=1}^N \lambda_i \gamma(\mathbf{h}(i, 0))$$

Note that replacing covariances by semivariances changes the sign.

Done! This is now a **computable** expression for the **kriging variance** at any point \mathbf{x}_0 , given the **locations** of the sample points \mathbf{x}_i , once the **weights** λ_i are known.

Next we will see how to use this expression to select optimum λ_i .

To check your understanding ...

Q4 : *Does the kriging variance depend on the **data values** at the sample points, or on the predicted data value? How can you see this from the equation?* *Jump to A4 •*

Summary: OK Kriging Variance

- Depends on the variogram function $\gamma(\mathbf{h})$ and the point configuration **around each point** to be predicted; we have derived this as:

$$\sigma^2(Z(\mathbf{x}_0)) \equiv E[\{\hat{Z}(\mathbf{x}_0) - Z(\mathbf{x}_0)\}^2] = 2 \sum_{i=1}^N \lambda_i \gamma(\mathbf{x}_i, \mathbf{x}_0) - \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j \gamma(\mathbf{x}_i, \mathbf{x}_j)$$

- * First term: **lower semi-variances between a point and the sample points leads to a lower kriging variance**; different for each point to be predicted
 - * Second term: **respect the co-variance structure of the sample points**; depends on configuration of sample points only
- We do not yet know what are the optimal weights λ , but once we do, we can calculate this kriging variance; so **we can selected the λ to minimize it**.

‘Model globally, predict locally’

- The **kriging equations** are solved **separately for each point** \mathbf{x}_0 , using the semivariances around that point, in a local neighbourhood; this gives a **different** set of weights λ for each point to be predicted.
- However, the **variogram model** $\gamma()$ used in these equations is estimated **only once**, using information about the spatial structure over the whole study area.
- Q: How is this possible?
 - * A1: Assume **second-order stationarity**
 - * A2: Assume at least **local first-order stationarity** (local weights will be high enough to mask long-distance non-stationarity)

Computing the weights

- There is one important piece of the puzzle missing: **How do we compute the weights** λ to predict at a given point?
 - * (Recall the *ad-hoc* method: some power of inverse distance)
- We want these to be the “best”, based on a computable **objective function**.
- There will be an optimum combination of weights at the point to be predicted, given the **point configuration** and the **modelled variogram**
- We compute these weights for each point to be predicted, by an **optimization criterion**, which in OK is **minimizing the kriging variance**.
- The previous slides have shown how to derive a **computable expression** for the kriging variance.

Commentary

The next several slides require a knowledge of **differential calculus**; however the central idea is simple: to **minimizing** an **objective function** (which we have just derived, i.e. the kriging variance).

The techniques of differential calculus are used to do the minimization.

Objective function (1): Unconstrained

In a minimization problem, we must define an **objective function** to be minimized. In this case, it is the kriging variance in terms of the N weights λ_i :

$$f(\lambda) = 2 \sum_{i=1}^N \lambda_i \gamma(\mathbf{x}_i, \mathbf{x}_0) - \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j \gamma(\mathbf{x}_i, \mathbf{x}_j)$$

But this is **unbounded** and can be trivially solved by setting all weights to 0. We must add **another constraint** to **bound** it.

Objective function (2): Constrained

To bound the objective function, we need another constraint; here it is naturally **unbiasedness**; that is, the **weights must sum to 1**.

This is added to the system with a **LaGrange multiplier** ψ :

$$f(\lambda, \psi) = 2 \sum_{i=1}^N \lambda_i \gamma(\mathbf{x}_i, \mathbf{x}_0) - \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j \gamma(\mathbf{x}_i, \mathbf{x}_j) - 2\psi \left\{ \sum_{i=1}^N \lambda_i - 1 \right\}$$

Note that the last term $\equiv 0$, i.e. the prediction is unbiased. The LaGrange multiplier may be changed (we will see how) but that final term always drops out of the prediction.

Note: It is possible to derive the kriging system in other ways, without the LaGrange multipliers and in terms of covariances. See Topic “Regression approach to kriging” below.

Minimization

This is now a system of $N + 1$ equations in $N + 1$ unknowns.

Minimize by setting all $N + 1$ **partial derivatives** to zero:

$$\frac{\partial f(\lambda_i, \psi)}{\partial \lambda_i} = 0, \forall i$$
$$\frac{\partial f(\lambda_i, \psi)}{\partial \psi} = 0$$

In the differential equation with respect to ψ , all the λ are constants, so the first two terms differentiate to 0; in the last term the ψ differentiates to 1 and we are left with the unbiasedness condition:

$$\sum_{i=1}^N \lambda_i = 1$$

The Kriging system

In addition to unbiasedness, the partial derivatives with respect to the λ_i give N equations (one for each λ_i) in $N + 1$ unknowns (the λ_i plus the LaGrange multiplier ψ):

$$\sum_{j=1}^N \lambda_j \gamma(\mathbf{x}_i, \mathbf{x}_j) + \psi = \gamma(\mathbf{x}_i, \mathbf{x}_0), \quad \forall i$$

This is now a **system of $N + 1$ equations in $N + 1$ unknowns** and can be solved by linear algebra.

The **semivariances between sample points** $\gamma(\mathbf{x}_i, \mathbf{x}_j)$ are computed **only once** for any point configuration; however the **semivariances at a sample point** $\gamma(\mathbf{x}_i, \mathbf{x}_0)$ must be **computed separately for each point to be predicted**.

Solving the Kriging system

At each point to be predicted:

1. **Compute the semivariances** γ from the separation between the point and the samples, according to the **modelled variogram**
2. **Solve simultaneously** for the weights and multiplier
3. **Compute the predicted value** as the **weighted average** of the samples
4. **Compute the variable term** of the kriging variance
5. **Add the constant term** of the kriging variance to get the total variance.

Importance of the variogram model

- The kriging system is solved using the modelled semi-variances
- **Different models will give different kriging weights** to the sample points ...
- ...and these will give different predictions
- Conclusion: **bad model → bad predictions**

Matrix form of the Ordinary Kriging system

$$\mathbf{A}\lambda = \mathbf{b}$$

$$\mathbf{A} = \begin{bmatrix} \gamma(\mathbf{x}_1, \mathbf{x}_1) & \gamma(\mathbf{x}_1, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_1, \mathbf{x}_N) & 1 \\ \gamma(\mathbf{x}_2, \mathbf{x}_1) & \gamma(\mathbf{x}_2, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_2, \mathbf{x}_N) & 1 \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ \gamma(\mathbf{x}_N, \mathbf{x}_1) & \gamma(\mathbf{x}_N, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_N, \mathbf{x}_N) & 1 \\ 1 & 1 & \cdots & 1 & 0 \end{bmatrix}$$

$$\lambda = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_N \\ \psi \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} \gamma(\mathbf{x}_1, \mathbf{x}_0) \\ \gamma(\mathbf{x}_2, \mathbf{x}_0) \\ \vdots \\ \gamma(\mathbf{x}_N, \mathbf{x}_0) \\ 1 \end{bmatrix}$$

Inside the Matrix

The **block matrix** notation shows the semivariances and LaGrange multiplier explicitly:

$$\mathbf{A} = \begin{pmatrix} \Gamma & \mathbf{1} \\ \mathbf{1}^T & 0 \end{pmatrix}$$

$$\lambda = \begin{bmatrix} \Lambda \\ \psi \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} \Gamma_0 \\ 1 \end{bmatrix}$$

Solution

This is a system of $N + 1$ equations in $N + 1$ unknowns, so can be solved uniquely, as long as \mathbf{A} is positive definite; this is guaranteed by using authorized models. This has the solution (in matrix notation):

$$\boldsymbol{\lambda} = \mathbf{A}^{-1}\mathbf{b}$$

Now we can **predict** at the point, using the weights:

$$\hat{Z}(\mathbf{x}_0) = \sum_{i=1}^N \lambda_i z(\mathbf{x}_i)$$

The **kriging variance** at a point is given by the scalar product of the weights (and multiplier) vector $\boldsymbol{\lambda}$ with the right-hand side of the kriging system: Note that $\boldsymbol{\lambda}$ includes as its last element ψ , which depends on covariance structure of the sample points:

$$\hat{\sigma}^2(\mathbf{x}_0) = \mathbf{b}^T \boldsymbol{\lambda}$$

Exercise

At this point you should do the **first section** of **Exercise 5: Predicting from point samples (Part 2)** which is provided on the module CD:

- §1 **Kriging weights**

This should take less than an hour.

As in all exercises there are **Tasks**, followed by R code on how to complete the task, then some **Questions** to test your understanding, and at the end of each section the **Answers**. Make sure you understand all of these.

Topic: Simple Kriging

Recall that In OK:

- We must estimate the **regional** (stationary) mean along with the predicted values, in the OK system.
- However, there may be situations where **the regional mean is known**. Then we can use so-called **Simple Kriging** (SK)

Similarly for UK or KED:

- We must estimate both the intercept (β_0) and all other trend coefficients (β_i), along with the predicted values, in the UK system.
- Similarly, if **the trend is known**, we can use “Simple” variants of UK and KED.

Simple Kriging (SK)

- The (stationary) **regional mean** must be known *a priori*
 - * in **Regression Kriging**, by definition the residuals have mean 0.
 - * in **Indicator Kriging** we expect each quantile to have the corresponding proportion of 1's (if the sample was unbiased)
 - * Note: the mean of a spatial sample is generally *not* the spatial mean, precisely because the observations are spatially-correlated.
- The mean may also be known from **previous studies** that give a better estimate than the (small) sample size we are working with; this could be a study of a larger area enclosing our study area.
- A **non-stationary** (moving) mean may be known from external evidence, e.g. a modelling exercise or regression analysis.

SK Predictions

- We reformulate the OK estimate **without the constraint** that weights sum to 1
- So, any bias from the weights must be compensated with respect to the known mean μ when predicting at a point:

$$\hat{Z}_{SK}(\mathbf{x}_0) = \sum_{i=1}^N \lambda_i z(\mathbf{x}_i) + \{1 - \sum_{i=1}^N \lambda_i\} \mu$$

- I.e. **Estimate = Weighted Sum + Bias Correction**
- The correction term is not present in OK prediction; its analog is the constraint that the weights sum to 1.

The SK system

There is **no need for a LaGrange multiplier** in the SK system, since there is no unbiasedness condition on the weights at any point. The system thus has one less row and column than the OK system.

$$\mathbf{A}\boldsymbol{\lambda} = \mathbf{b}$$

$$\mathbf{A} = \begin{bmatrix} \gamma(\mathbf{x}_1, \mathbf{x}_1) & \gamma(\mathbf{x}_1, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_1, \mathbf{x}_N) \\ \gamma(\mathbf{x}_2, \mathbf{x}_1) & \gamma(\mathbf{x}_2, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_2, \mathbf{x}_N) \\ \gamma(\mathbf{x}_N, \mathbf{x}_1) & \gamma(\mathbf{x}_N, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

$$\boldsymbol{\lambda} = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_N \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} \gamma(\mathbf{x}_1, \mathbf{x}_0) \\ \gamma(\mathbf{x}_2, \mathbf{x}_0) \\ \vdots \\ \gamma(\mathbf{x}_N, \mathbf{x}_0) \end{bmatrix}$$

Solution of the SK system

This is a system of N equations in N unknowns, so can be solved uniquely, as long as \mathbf{A} is positive definite (as in OK); this is guaranteed by using authorized models.

However, since the weights λ_i are not constrained to sum to 1, in addition **the variogram must be bounded** (e.g. an exponential or spherical model; but not a power model).

Then:

$$\boldsymbol{\lambda} = \mathbf{A}^{-1}\mathbf{b}$$

Estimate at the point, using the weights and the known mean:

$$\hat{Z}_{SK}(\mathbf{x}_0) = \sum_{i=1}^N \lambda_i z(\mathbf{x}_i) + \left\{1 - \sum_{i=1}^N \lambda_i\right\} \mu$$

Note that the unbiasedness at each prediction point is restored here; if the λ did not sum to 1, it is corrected.

SK variance

The kriging variance at a point is given by:

$$\begin{aligned}\hat{\sigma}_{SK}^2(\mathbf{x}_0) &= \mathbf{b}^T \boldsymbol{\lambda} \\ &= \sum_{i=1}^N \lambda_i \gamma(\mathbf{x}_i, \mathbf{x}_0)\end{aligned}$$

This formula has a simple interpretation:

- the kriging variance is the **weighted sum** of the **semivariances** between each sample point and the prediction point;
- the weights are the **kriging weights** found by solving the kriging system.

Topic: Block Kriging

Often we want to predict **average values** of some target variable in **blocks** of some defined size, not at points.

Example: average woody biomass in a forest block of 40ha, if this is a minimum management unit, e.g. we will decide to harvest or not the whole block. We don't care about any finer-scale information, we wouldn't use it if we had it.

Block kriging (BK) is quite similar in form to OK, but the **kriging variances are lower**, because the **within-block variability** is removed.

Commentary

There is only one new idea in this section: by predicting an **average over some block** larger than the support, we **reduce the kriging variance**.

Most of this section shows how this is expressed mathematically.

The practical implication is shown in the comparative figures at the end of the section.

Block Ordinary Kriging (BK)

- Estimate at blocks of a defined size, with unknown mean (which must also be estimated) and no trend
- Each **block** B is estimated as the weighted average of the values at all sample **points** \mathbf{x}_i :

$$\hat{Z}(B) = \sum_{i=1}^N \lambda_i z(\mathbf{x}_i)$$

- As with OK, the weights λ_i sum to 1, so that the estimator is **unbiased**, as for OK

The Block Kriging system

The same derivation as for the OK system produces these equations:

$$\sum_{j=1}^N \lambda_j \gamma(\mathbf{x}_i, \mathbf{x}_j) + \psi(B) = \overline{\gamma}(\mathbf{x}_i, B), \quad \forall i$$

- The semivariances in the right-hand side, i.e., the **b** vector, are now between **sample points** and the **block** to be predicted
- The semivariance with a block is written as $\overline{\gamma}$, the overline indicating an **average**
- The left-hand side, i.e., the **A** matrix (semivariances between known observation *points*), is the same as in OK

Matrix form of the Block Kriging system

$$\mathbf{A}\lambda = \mathbf{b}$$

$$\mathbf{A} = \begin{bmatrix} \gamma(\mathbf{x}_1, \mathbf{x}_1) & \gamma(\mathbf{x}_1, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_1, \mathbf{x}_N) & 1 \\ \gamma(\mathbf{x}_2, \mathbf{x}_1) & \gamma(\mathbf{x}_2, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_2, \mathbf{x}_N) & 1 \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ \gamma(\mathbf{x}_N, \mathbf{x}_1) & \gamma(\mathbf{x}_N, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_N, \mathbf{x}_N) & 1 \\ 1 & 1 & \cdots & 1 & 0 \end{bmatrix}$$

$$\lambda = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_N \\ \psi(B) \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} \bar{\gamma}(\mathbf{x}_1, B) \\ \bar{\gamma}(\mathbf{x}_2, B) \\ \vdots \\ \bar{\gamma}(\mathbf{x}_N, B) \\ 1 \end{bmatrix}$$

Solution

$$\lambda = \mathbf{A}^{-1}\mathbf{b}$$

Now we can estimate at the block, using the weights:

$$\hat{Z}(B) = \sum_{i=1}^N \lambda_i z(\mathbf{x}_i)$$

The kriging variance for the block is given by:

$$\hat{\sigma}^2(B) = \mathbf{b}^T \lambda - \bar{y}(B, B)$$

Note that **the variance is reduced by the within-block variance.**

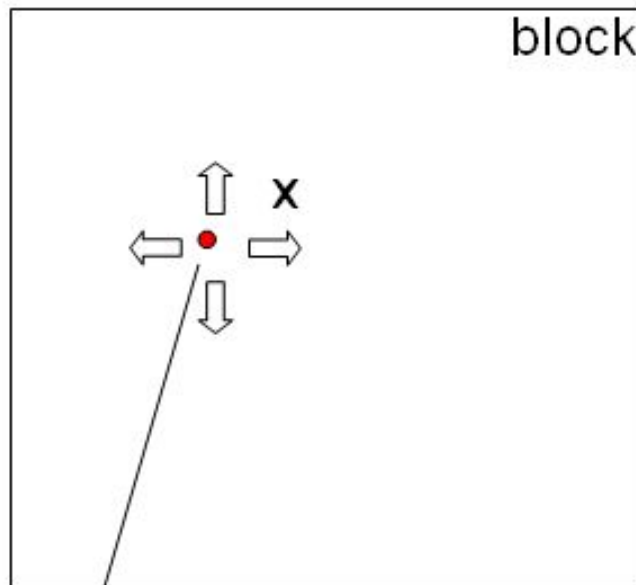
Commentary

The semivariances γ in the above formulation are not just a function of separation, because they are *not* between points. Instead, they are between sample **points** and prediction **blocks**. This is illustrated in the following figure (left side).

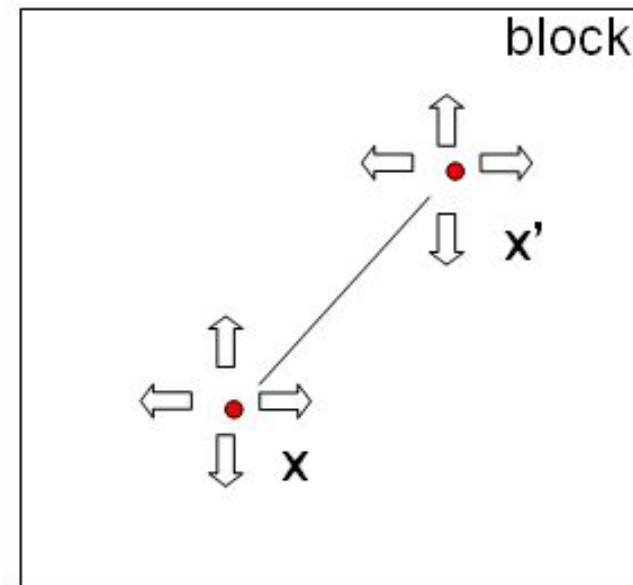
In addition, the variance at the prediction location is now not at a point, but rather at a **block**. So some of the kriging variance must be accounted for within that block, i.e. the variance that is due to short-range variability at distances **shorter than the block size**. This is illustrated in the following figure (right side).

We will then discuss in detail how to compute these two.

Integration of semivariances for block kriging



x_i integration
from sampling
point to block



integration within block

Computing the semivariance between a point and a block

The complication here, compared to OK, is that the semivariances in the **b** vector are between sample points and the entire **block** to be predicted: $\bar{\gamma}(\mathbf{x}_i, B)$

(Note that in the **A** matrix the semivariances are between known sample *points*, so the point-to-point semivariance $\gamma(\mathbf{x}_i, \mathbf{x}_j)$ is used, as in the OK system.)

So there is not a single distance than can be substituted into the variogram model. We have to integrate over the block:

$$\bar{\gamma}(\mathbf{x}_i, B) = \frac{1}{|B|} \int_B \gamma(\mathbf{x}_i, \mathbf{x}) d\mathbf{x}$$

where $|B|$ is the area of the block, and \mathbf{x} is a point within the block.

As written **all** points in the block (conceptually, an infinite number!) would be used; in practice, this is done by **discretization** of the block into a set of points.

Computing the within-block variances

This is the factor by which the estimation variance is reduced:

$$\bar{\gamma}(B, B) = \frac{1}{|B|^2} \int_B \int_B \gamma(\mathbf{x}, \mathbf{x}') d\mathbf{x} d\mathbf{x}'$$

As the block size $|B|$ approaches zero, the double integral also approaches zero; in fact this is the limit. This shows that **OK is a special case of BK**.

In practice, this is also calculated by **discretizing** the block into n points:

$$\bar{\gamma}(B, B) \approx \frac{1}{|B|^2} \sum_{i=1}^n \sum_{j=1}^n w_i w_j \gamma(\mathbf{x}_i, \mathbf{x}_j)$$

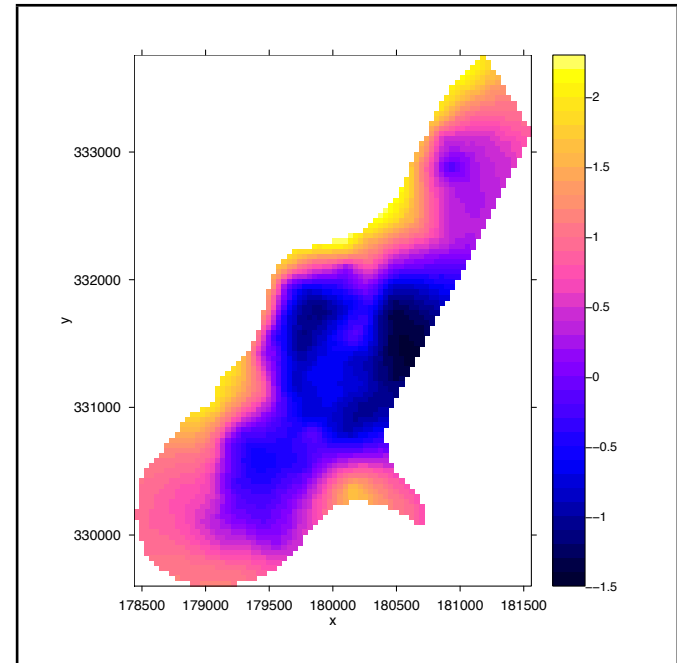
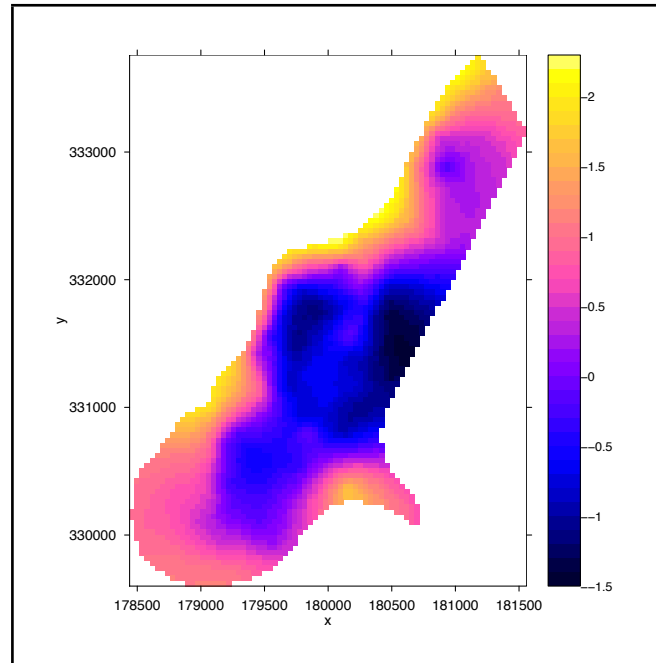
where of course $\sum_i w_i = 1$; the weights are set by their position within the unit block.

Visualizing the effect of block size

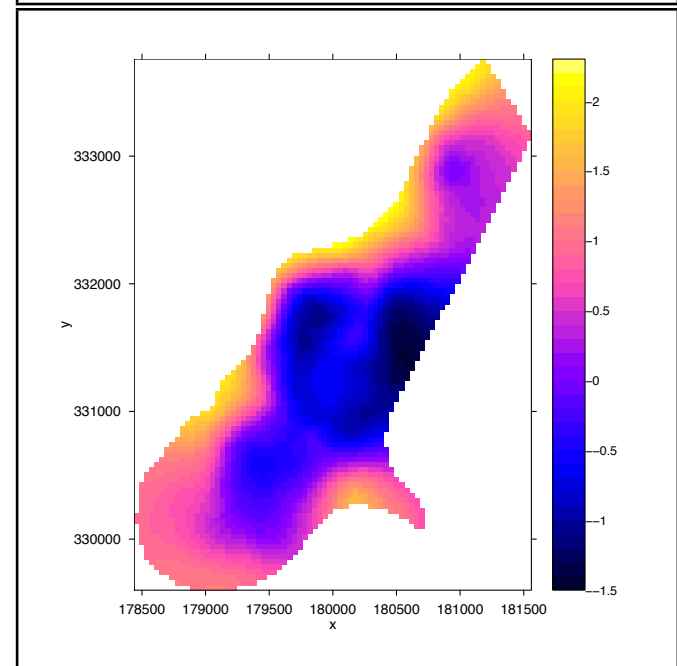
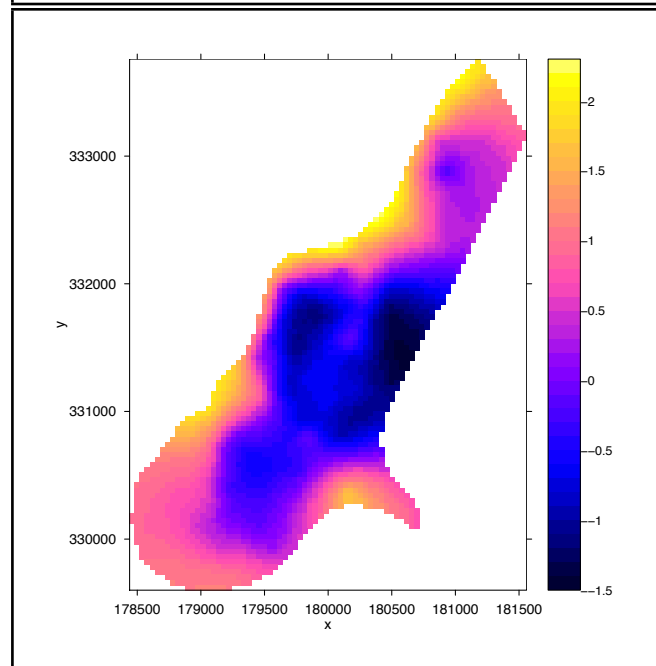
The following graphs show the changing **predictions** and their **variances** as the block size is increased. This is from the Meuse soil pollution study; target variable is $\log_{10}(\text{Cd})$.

Each graph uses one scale, to allow direct comparison.

Predictions: OK, BK10



BK40, BK160

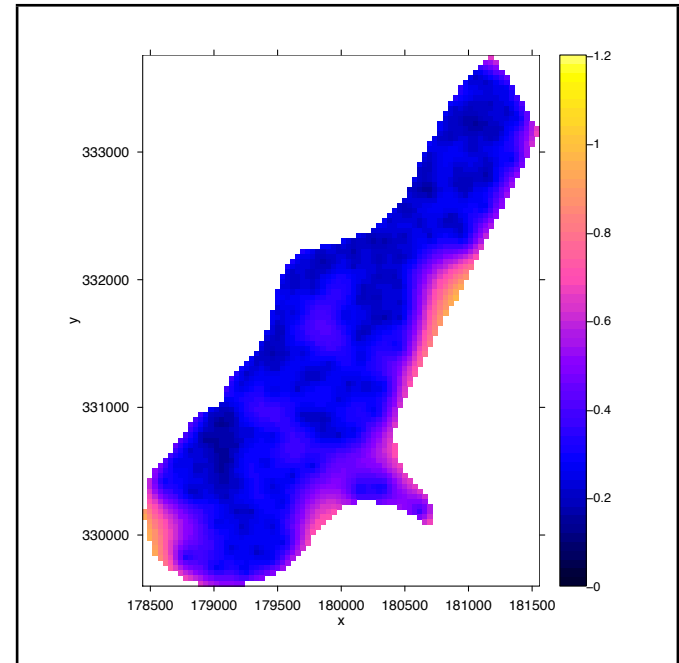
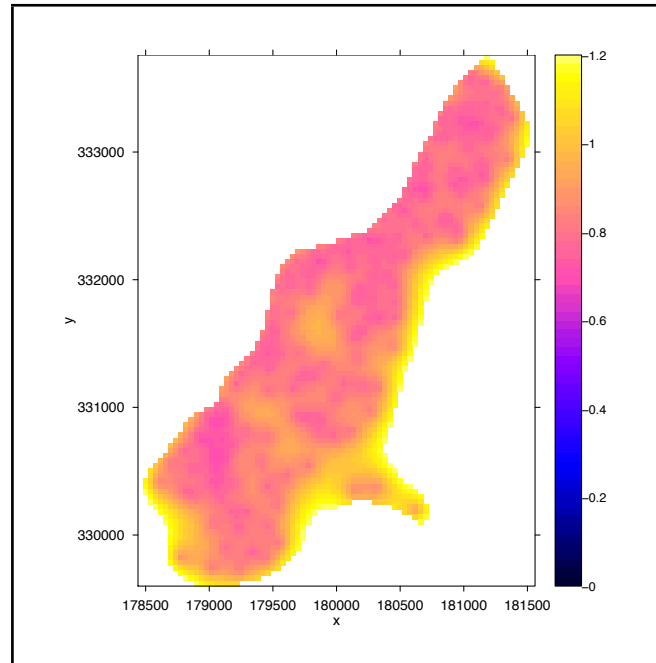


To check your understanding ...

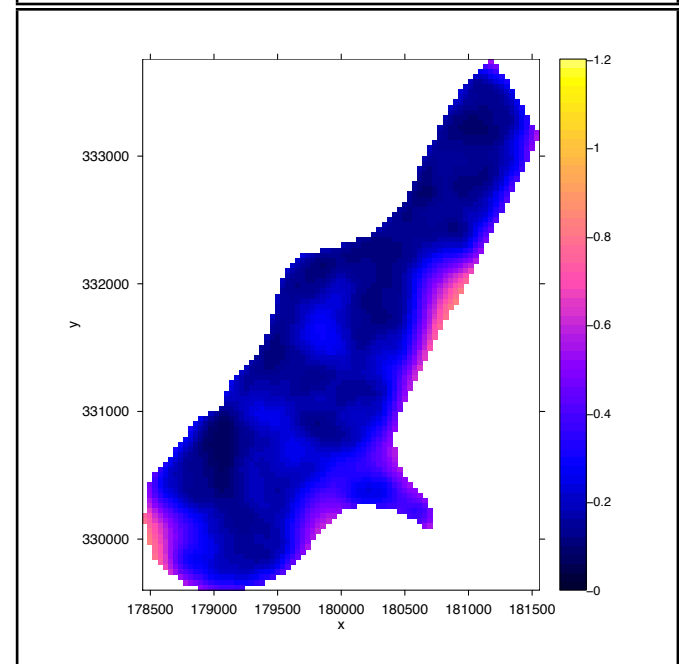
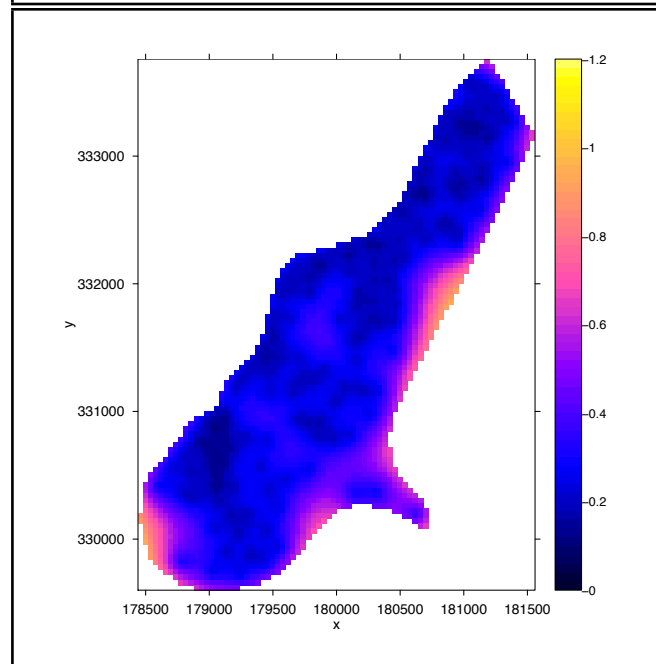
Q5 : *What are the differences between the **predictions** with different block sizes?*

Jump to A5 •

Variances: OK, BK10



BK40, BK160



To check your understanding ...

Q6 : *What are the differences between the **kriging variances** with different block sizes? Jump to A6 •*

Topic: Universal Kriging (UK)

This is a **mixed predictor** which includes a **global trend** as a function of the **geographic coördinates** in the kriging system, as well as **local spatial dependence**.

Example: The depth to the top of a given sedimentary layer may have a regional trend, expressed by geologists as the **dip** (angle) and **strike** (azimuth). However, the layer may also be locally thicker or thinner, or deformed, with spatial covariance in this local structure.

UK is recommended when there is evidence of **first-order non-stationarity**, i.e. the **expected value** varies across the map, but there is still **second-order stationarity** of the **residuals** from this trend.

Note: The “global” trend can also be fitted **locally**, within some user-defined radius, so that this interpolator can range from local (immediate neighbourhood) to global (whole area), according to the analyst’s evidence on spatial structure.

Terminology: UK vs. KED vs. RK

- In these notes, we restrict the term “Universal Kriging” (UK) to the use of a **geographic trend** as co-variate
- The term “UK” is used by some authors also to include a trend in one or more **feature-space** predictors, i.e. co-variables.
- In this course we call this use of feature-space predictors **Kriging with External Drift** (KED)
- **The mathematics are the same**; it is the co-variables that are different:
 - * UK : only geographic coördinates;
 - * KED : also (or only) feature-space co-variables; solve trend and local structure in one kriging system;
 - * RK (Regression Kriging): may also solve trend and local structure separately.

An abstract view of UK/KED/RK

The realization of variable Z at spatial location \mathbf{s} can be considered as the result of three distinct processes:

$$Z(\mathbf{s}) = m(\mathbf{s}) + \epsilon'(\mathbf{s}) + \epsilon''(\mathbf{s})$$

m a **deterministic** component; e.g., a regional trend, or the effect of some forcing variable;

ϵ' a **spatially-correlated** stochastic process;

ϵ'' pure **noise**, not spatially-correlated, not deterministic

Note that all three processes are distributed in space, i.e., vary by location.

Relation of this formulation to trends and OK

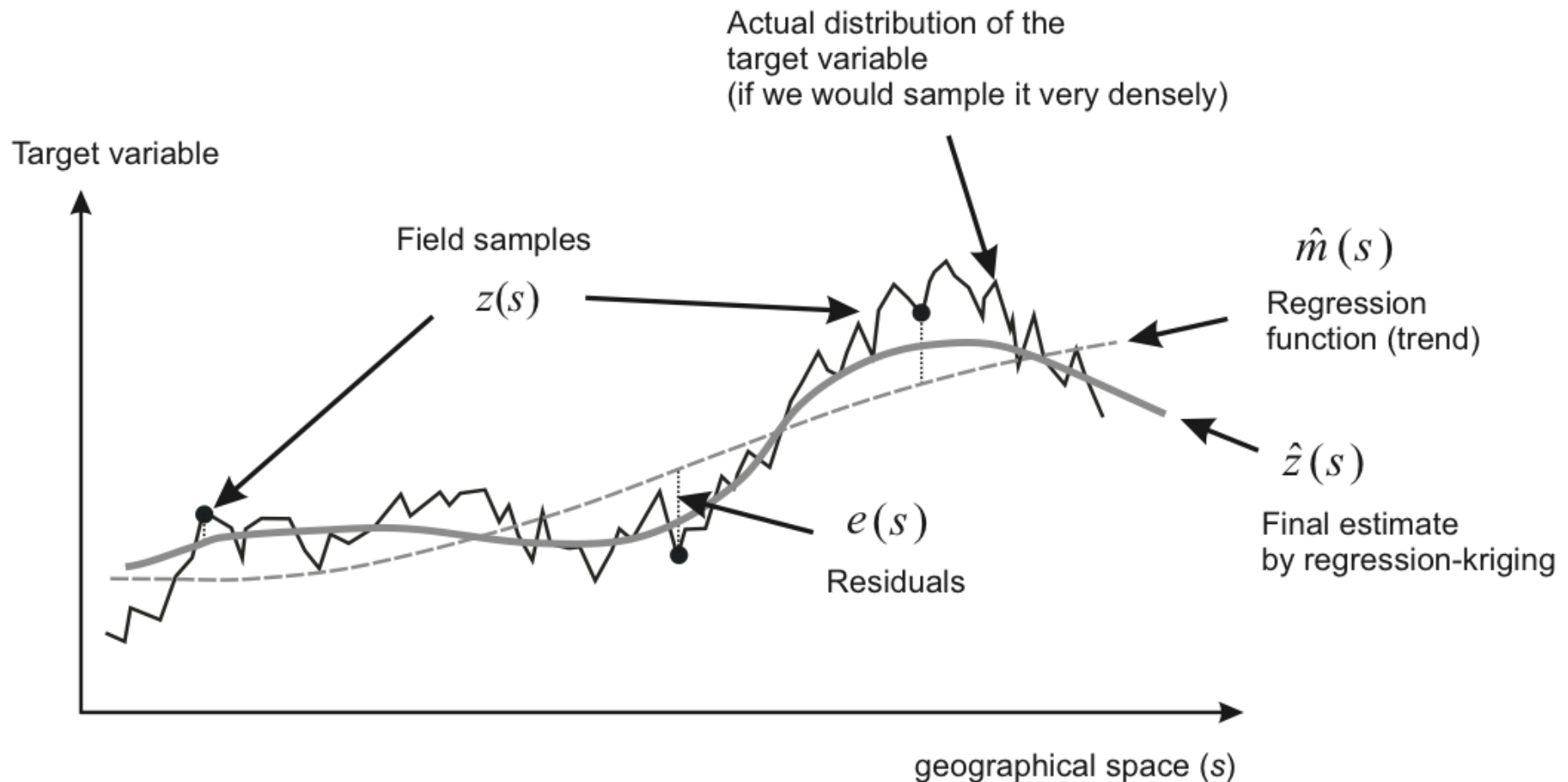
OK:

- $Z(\mathbf{s}) = \mu + \epsilon'(\mathbf{s}) + \epsilon''(\mathbf{s})$
- i.e., the deterministic part $m(\mathbf{s})$ is just a single expected value μ , the overall level of the target variable, with no spatial structure
- The noise term $\epsilon''(\mathbf{s})$ is the nugget variance

Trend surface:

- $Z(\mathbf{s}) = m(\mathbf{s}) + \epsilon''(\mathbf{s})$
- i.e., a deterministic trend (modelled from the coordinates) and noise
- The noise term $\epsilon''(\mathbf{s})$ is the lack of fit

Visualizing the abstract model



Source: Hengl, T. (2009). *A Practical Guide to Geostatistical Mapping*. Amsterdam.
Retrieved from <http://spatial-analyst.net/book/>; Figure 2.1

Note this figure uses the term “regression kriging” for the abstract model.

Prediction with UK

In UK, we model the value of variable z at location \mathbf{x}_i as the sum of:

- a regional **non-stationary trend** $m(\mathbf{x}_i)$ and a
- a local **spatially-correlated random component** $e(\mathbf{x}_i)$; the **residuals** from the regional trend

$$Z(\mathbf{x}_i) = m(\mathbf{x}_i) + e(\mathbf{x}_i)$$

Note that the random component is now expected to be **first-order stationary**, because non-stationarity is all due to the trend

Here $m(\mathbf{x})$ is **not a constant** as in OK, but instead is a **function of position**, i.e. the global **trend**.

Base functions

The trend is modelled as a **linear combination** of p known **base functions** $f_j(s)$ and p unknown constants β_j (these are the **parameters** of the base functions):

$$Z(\mathbf{x}_i) = \sum_{j=1}^p \beta_j f_j(\mathbf{x}_i) + e(\mathbf{x}_i)$$

Examples of base functions

- For **linear** drift:

$$f_0(\mathbf{x}) = 1, f_1(\mathbf{x}) = x_1, f_2(\mathbf{x}) = x_2$$

where x_1 is one coördinate (say, E) and x_2 the other (say, N)

- Note that $f_0(\mathbf{x}) = 1$ estimates the global mean (as in OK).
- For **quadratic** drift: also include second-order terms:

$$f_3(\mathbf{x}) = x_1^2, f_4(\mathbf{x}) = x_1x_2, f_5(\mathbf{x}) = x_2^2$$

Predictions at points

A point is **predicted** as in OK:

$$\hat{Z}(\mathbf{x}_0) = \sum_{i=1}^N \lambda_i z(\mathbf{x}_i)$$

But, **the weights** λ_i for each sample point take into account both the **global trend** and **local effects**.

We need to set the UK system up to include both of these.

Unbiasedness of predictions

The **unbiasedness** condition is expressed with respect to the **trend** as well as the overall mean (as in OK):

$$\sum_{i=1}^N \lambda_i f_k(\mathbf{x}_i) = f_k(\mathbf{x}_0), \quad \forall k$$

The expected value at each point of all the **functions** must be that predicted by that function. The first of these is the overall mean (as in OK).

Example: If $f_1(\mathbf{x}_0) = x_1$ (linear trend towards the E), then at each point \mathbf{x}_0 the expected value must be x_1 , i.e. the point's E coördinate:

$$\sum_{i=1}^N \lambda_i x_i = x_1$$

This is a **further restriction** on the weights λ .

Computing the experimental semivariogram for UK

- The semivariances γ are based on the **residuals**, not the original data, because the **random field** part of the spatial structure applies only **after** any trend has been removed.
- How to obtain?
 1. Calculate the **best-fit surface**, with the same base functions to be used in UK;
 2. **Subtract** the trend surface at the data points from the data value to get residuals;
 3. Compute the **variogram** of the **residuals**.

Note: Some programs, e.g. gstat, will do these all in one step.

Characteristics of the residual variogram

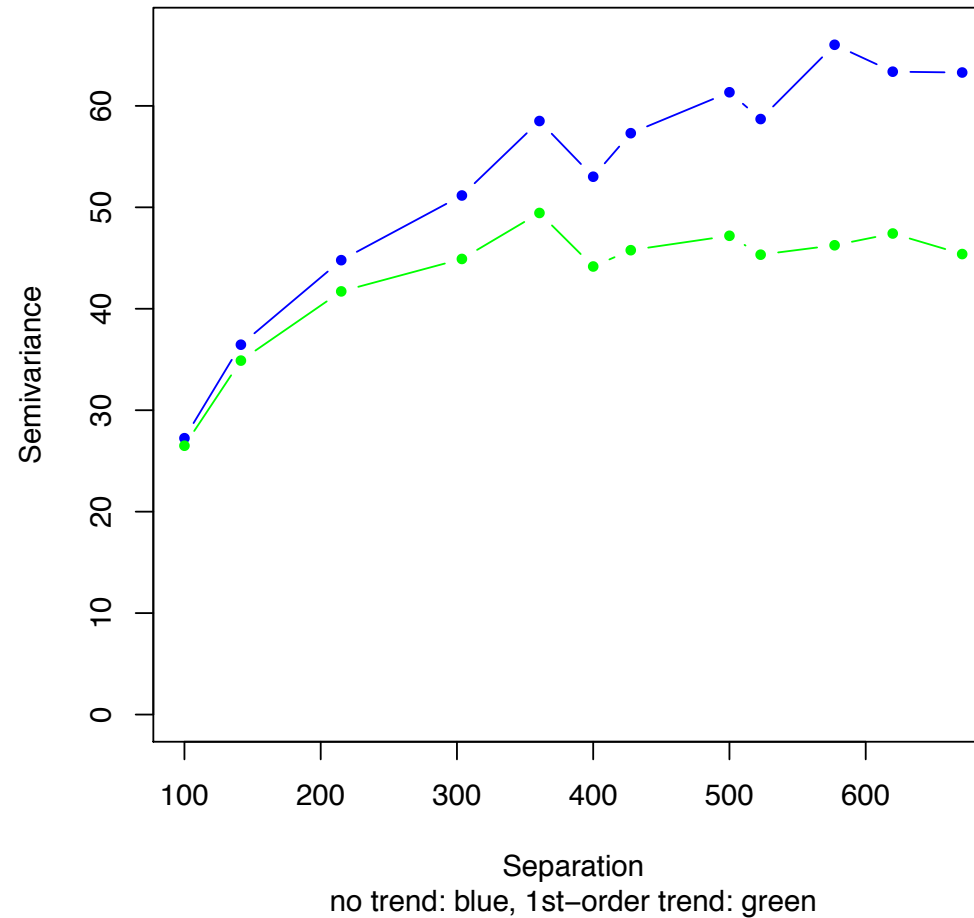
If there is a strong trend, the variogram model **parameters** for the residuals will be very different from the original variogram model, since **the global trend has taken out some of the variation**, i.e. that due to the long-range structure.

The usual case is:

- lower **sill** (less total variability)
- shorter **range** (long-range structure removed)

Example original vs. residual variogram

Variograms, Oxford soils, CEC (cmol+ kg⁻¹ soil)



To check your understanding . . .

Q7 : *What are the approximate sill and range of the original (blue) and residual (green) variograms?*

The first-order trend surface in this example explained about 40% of the overall variance in the target variable. How is this reflected in the variogram of the residuals from this surface (as shown)? Jump to A7 •

Q8 : *What is the relation between the nugget variance of the original (blue) and residual (green) variograms? What should the relation be, in theory?* Jump to A8 •

Universal Kriging: Local vs. Global trends

As with OK, UK can be used two ways:

- **Globally**: using **all** sample points when predicting each point
- **Locally**, or in **patches**: restricting the sample points used for prediction to some **search radius** (or sometimes **number of neighbours**) around the point to be predicted

We now discuss the properties of each:

UK over the region

Using **all** sample points when predicting each point:

- Appropriate if there is a **regional trend** across the entire study area
- Agrees with the global computation of the residual variogram
- This gives the same results as Regression Kriging on the coördinates

UK in a neighbourhood

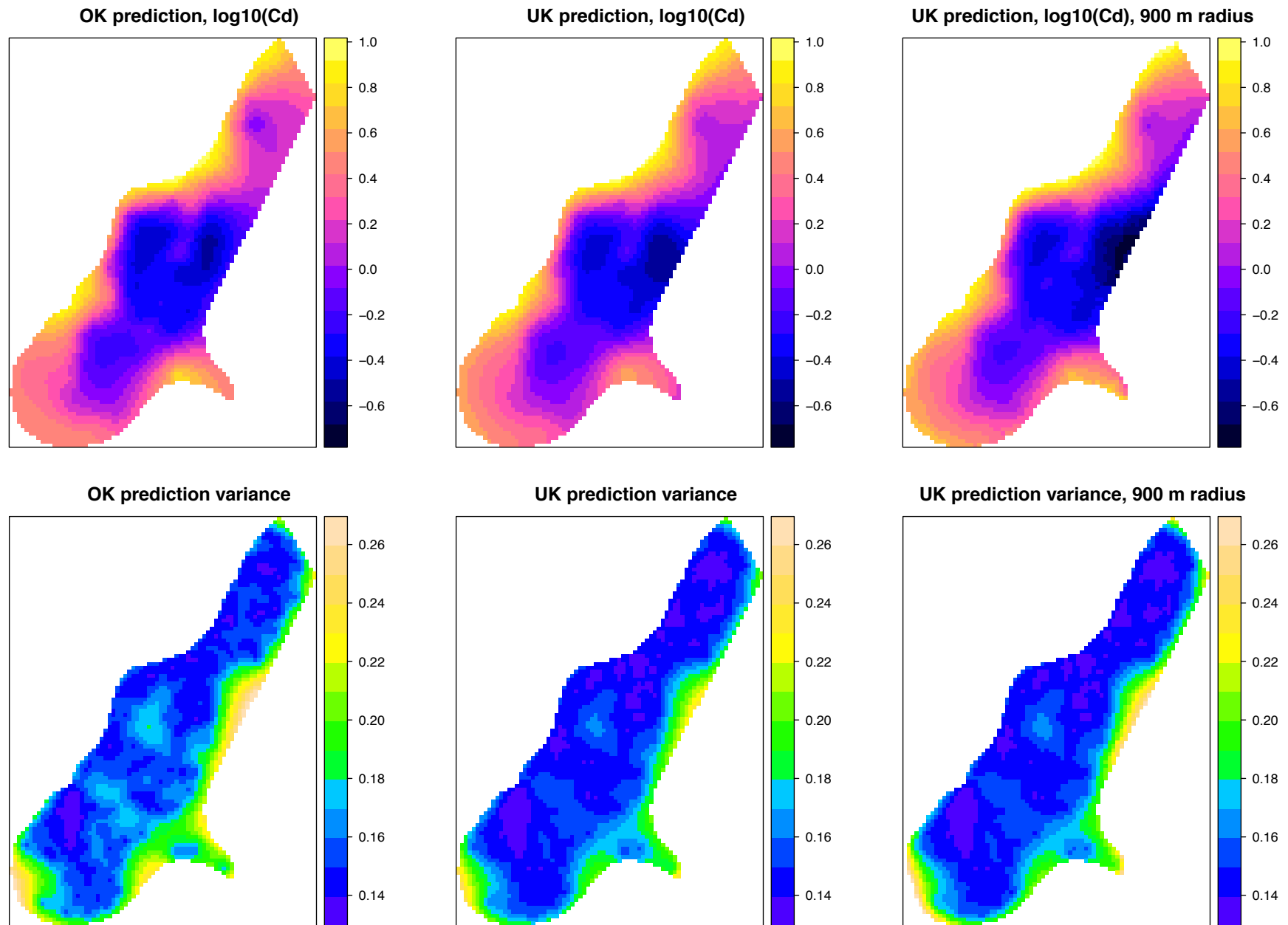
Using just the points in some **neighbourhood**:

- This allows the **trend surface to vary** over the study area, since it is **re-computed at each prediction point**
- Appropriate to smooth away some **local variation in a trend**, difficult to justify theoretically
- Note that the **residual variogram** was not computed in patches, but assuming a global trend
- Leads to some patchiness in the map
- There should be some **evidence of patch size**, perhaps from the original (*not* residual) variogram; this can be used as the search radius.

Compare UK to OK on the next page; the NW-SE trend somewhat modifies the predictions. Note the lower kriging variances with UK, due to the trend surface.

UK with a local trend is intermediate between global UK and OK.

All predictions are shown with the same scale, similarly for variances.



To check your understanding . . .

Q9 : *What are the major differences in the above figure between OK, global UK, and neighbourhood UK predictions?*

Jump to A9 •

Commentary

UK in a neighbourhood is more difficult to justify than OK in a neighbourhood.

In the case of OK, we have a variogram showing the effective **range** of the local spatial dependence, so that points further than this from the prediction point only contribute to the expected value.

In the case of UK, the residual variogram is computed from the **global**, not local, trend. The range is generally much shorter than the “effective trend” near a prediction point.

Topic: Derivation of the Universal Kriging (UK) system

Again, two approaches:

1. Regression (linear algebra)
2. Minimization (differential calculus)

These are generalizations of the derivations of the OK system, above.

Regression approach to UK

The general linear model applied to OK can be extended to UK (or mathematically-equivalent KED). Instead of a vector of 1's we have the **design matrix** \mathbf{Q} and the values \mathbf{x}_0 of the **basis functions** at the prediction point.

The spatial mean $\hat{\mu}$ is replaced by all the parameters $\hat{\beta}$ of the trend:

$$\hat{\beta} = (\mathbf{Q}^T \cdot \mathbf{C}^{-1} \cdot \mathbf{Q})^{-1} \cdot (\mathbf{Q}^T \cdot \mathbf{C}^{-1} \cdot \mathbf{z})$$

Note that the **GLS trend surface coefficients** are indeed produced by this method; analogous to the spatial mean in the regression formulation of OK.

This is an advantage of this formulation compared to the minimization approach.

(continued ...)

UK prediction (regression formulation)

Then the kriging prediction becomes:

$$\hat{Z}_{OK} = \mathbf{x}_0^T \hat{\boldsymbol{\beta}} + \mathbf{c}_0^T \cdot \mathbf{C}^{-1} \cdot (\mathbf{z} - \mathbf{Q}\hat{\boldsymbol{\beta}})$$

- $\mathbf{x}_0^T \hat{\boldsymbol{\beta}}$ here replaces the single estimate of the spatial mean $\hat{\mu}$ in OK; this term is the prediction from the GLS trend at the point.
- $(\mathbf{z} - \mathbf{Q}\hat{\boldsymbol{\beta}})$ here replaces $(\mathbf{z} - \hat{\mu}\mathbf{1})$ in OK; this is again the residual from the GLS model.

UK weights (regression formulation)

From this we derive the vector of kriging **weights**, i.e. what to multiply each observation by in the weighted sum (prediction):

$$\lambda^T = \mathbf{c}_0^T \cdot \mathbf{C}^{-1} - \mathbf{x}_0^T \cdot (\mathbf{Q}^T \cdot \mathbf{C}^{-1} \cdot \mathbf{Q})^{-1} \mathbf{Q}^T \cdot \mathbf{C}^{-1}$$

Note: derive this by collecting all the terms that multiply the observed values \mathbf{z} in the OK prediction.

UK prediction variance (regression formulation)

And similarly for the kriging variance:

$$\text{Var}(\hat{Z}_0 - Z_0) = c_{00} - \mathbf{c}_0^T \cdot \mathbf{C}^{-1} \cdot \mathbf{c}_0 + \frac{\mathbf{x}_0 - (\mathbf{c}_0^T \cdot \mathbf{C}^{-1} \cdot \mathbf{Q})}{\mathbf{Q}^T \cdot \mathbf{C}^{-1} \cdot \mathbf{Q}}$$

Note that this variance includes that due to the global trend; we can see this from the presence of the \mathbf{Q} (design) matrix, which includes the coordinates of each point.

For **block kriging**, replace c_{00} with c_{BB} , the average within-block variance (which will in general be smaller than the at-point variance).

Minimization approach to UK

As for OK, there is also a **minimization** approach to formulating the UK system.

The **objective function** to be **minimized** is the **kriging variance**; in addition to the unbiasedness constraint on the mean, there are unbiasedness constraints on the values of the base functions:

$$\sum_{j=1}^N \lambda_j \gamma(\mathbf{x}_i, \mathbf{x}_j) + \psi_0 + \sum_{k=1}^K \psi_k f_k(\mathbf{x}_i) = \gamma(\mathbf{x}_0, \mathbf{x}_i), \forall i \quad (\text{sample points}) \quad (1)$$

$$\sum_{i=1}^N \lambda_i = 1 \quad (\text{unbiasedness of mean}) \quad (2)$$

$$\sum_{i=j}^N \lambda_j f_k(\mathbf{x}_j) = f_k(\mathbf{x}_0), \forall k \quad (\text{base functions}) \quad (3)$$

Where is the trend surface?

Note that the coefficients β of the global trend surface are **not** mentioned in the UK system.

However, they are **implicit** in the solution and will affect the λ_i at each prediction point; that is, the local weights around each point will take into account the global trend.

(This has the advantage that we can restrict the search neighbourhood and fit the “global” trend only in a user-defined neighbourhood, i.e. in **patches**.)

The UK system (minimization formulation) (1)

We now present the mathematics of the UK system; it is not necessary to fully understand this in order to apply UK correctly.

The discussion uses similar notation to that for the OK system in the previous lecture.

Solve: $\mathbf{A}_U \lambda_U = \mathbf{b}_U$, where

$$\mathbf{A}_U = \begin{bmatrix} \gamma(\mathbf{x}_1, \mathbf{x}_1) & \cdots & \gamma(\mathbf{x}_1, \mathbf{x}_N) & 1 & f_1(\mathbf{x}_1) & \cdots & f_k(\mathbf{x}_1) \\ \vdots & \cdots & \vdots & \vdots & \vdots & \cdots & \vdots \\ \gamma(\mathbf{x}_N, \mathbf{x}_1) & \cdots & \gamma(\mathbf{x}_N, \mathbf{x}_N) & 1 & f_1(\mathbf{x}_N) & \cdots & f_k(\mathbf{x}_N) \\ 1 & \cdots & 1 & 0 & 0 & \cdots & 0 \\ f_1(\mathbf{x}_1) & \cdots & f_1(\mathbf{x}_N) & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ f_k(\mathbf{x}_1) & \cdots & f_k(\mathbf{x}_N) & 0 & 0 & \cdots & 0 \end{bmatrix}$$

The upper-left block $N \times N$ block is the spatial correlation structure (as in OK); the lower-left $k \times n$ block (and its transpose in the upper-right) are the trend predictor values at sample points; the rest of the matrix fits with λ_U and \mathbf{b}_U to set up the solution.

The UK system (minimization formulation) (2)

$$\lambda_U = \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_N \\ \psi_0 \\ \psi_1 \\ \vdots \\ \psi_k \end{bmatrix} \quad \mathbf{b}_U = \begin{bmatrix} y(\mathbf{x}_1, \mathbf{x}_0) \\ \vdots \\ y(\mathbf{x}_N, \mathbf{x}_0) \\ 1 \\ f_1(\mathbf{x}_0) \\ \vdots \\ f_k(\mathbf{x}_0) \end{bmatrix}$$

The λ_U vector contains the N weights for the sample points and the $k + 1$ LaGrange multipliers (1 for the overall mean and k for the trend model), and \mathbf{b}_U is structured like an additional column of \mathbf{A}_u , but referring to the point to be predicted.

UK system as an extension of OK (minimization formulation):

(1) Semivariance matrix

$$\mathbf{A_U} = \begin{pmatrix} \Gamma & \mathbf{1} & \mathbf{F} \\ \mathbf{1}^T & 0 & \mathbf{0}^T \\ \mathbf{F}^T & \mathbf{0} & \mathbf{0} \end{pmatrix}$$

Only the upper-left is used in OK:

$$\mathbf{A_O} = \begin{pmatrix} \Gamma & \mathbf{1} \\ \mathbf{1}^T & 0 \end{pmatrix}$$

Dimensions:

$$(n + 1 + k)^2 = \begin{pmatrix} n \times n & n \times 1 & n \times k \\ 1 \times n & 1 \times 1 & 1 \times k \\ k \times n & k \times 1 & k \times k \end{pmatrix}$$

UK system as an extension of OK (minimization formulation): (2) Weights and RHS

$$\lambda_U = \begin{bmatrix} \Lambda \\ \psi_0 \\ \Psi \end{bmatrix} \quad \lambda_O = \begin{bmatrix} \Lambda \\ \psi_0 \end{bmatrix}$$

$$\mathbf{b}_U = \begin{bmatrix} \Gamma_0 \\ 1 \\ \mathbf{F}_0 \end{bmatrix} \quad \mathbf{b}_O = \begin{bmatrix} \Gamma_0 \\ 1 \end{bmatrix}$$

UK (minimization formulation): Solution & kriging variance

Exactly as for OK:

$$\lambda_U = \mathbf{A}_U^{-1} \mathbf{b}_U$$

$$\sigma_U^2 = \mathbf{b}_U^T \lambda$$

This gives the Lagrange multipliers ψ for the overall mean and the trend, but not the trend surface coefficients β as they would be computed by the linear model. The prediction variance *does* include uncertainty in the trend.

Exercise

At this point you should do the **last sections** of **Exercise 5: Predicting from point samples (Part 2)** which is provided on the module CD:

- §2 **Block kriging**
- §3 **Universal kriging**

This should take about an hour.

As in all exercises there are **Tasks**, followed by R code on how to complete the task, then some **Questions** to test your understanding, and at the end of each section the **Answers**. Make sure you understand all of these.

Then do the **self-test** at the end of Exercise 5.

Topic: Kriging transformed variables

It may be desirable to **transform** a variable prior to variogram analysis and kriging.

Why transform?

- Ordinary Kriging assumes that the **distribution of deviations** from the single **expected value** is **Gaussian**
 - * The theory depends second-order stationary Gaussian **random fields**.
 - * Otherwise the kriging predictions (expected values) are unbiased but the **prediction variance** is not correct
 - * Confidence intervals can not be computed.
- OK is a **weighted average** of nearby values. A distribution that is not **symmetric** may present problems:
 - * In a highly **positively-skewed** distribution (often found in earth sciences) the few **high** values will overwhelm the others and lead to over-predictions.
- In a **multi-modal** distribution, the weighted average will predict intermediate values that fit none of the modes.

Thus, non-normal irregular distributions are often **transformed** to (approximate) normality.

Transformations

The most commonly-applied transformations are:

- **Logarithmic** for highly-skewed unimodal distributions
 - * Only if all values > 0 ; can add an offset before transforming
 - * For non-negative variables with some 0's, quite common to add half the detection/measurement limit
- **Square root** for 0 left-limited variables
- **Box-Cox**: logarithm and any power are special cases
- **Normal-score**, **rank-order**, or **Hermite polynomials** to convert irregular histograms to a normal distribution

Some attributes are already reported in transformed units – for example pH.

Lognormal kriging

1. The target variable is transformed as the natural logarithm: $y(\mathbf{x})_i = \log z(\mathbf{x}_i)$; this should be then approximately symmetrically distributed.
2. **Model** and **interpolate** with the transformed variable (OK, block kriging, UK, KED, ...); at each point we get a predicted values $\hat{Y}(\mathbf{x}_0)$ instead of $\hat{Z}(\mathbf{x}_0)$
3. **Back-transform** to original units.

The last step is optional if the transformed estimates can be used directly. For example, a regulatory threshold may be expressed directly in transformed units.

Back transformation

For the kriged estimate $\hat{Y}(\mathbf{x}_0)$, we want the back-transformed estimate $\hat{Z}(\mathbf{x}_0)$. We can not just exponentiate the estimate, because the estimate is a weighted **sum** (not product!) of logarithms. Further, the back-transformation (exponentiation) of a symmetric (Gaussian) prediction variance is right-skewed.

The back transformation is considered for two cases:

1. Simple Kriging (i.e. mean was **known a priori**)
2. Ordinary Kriging (i.e. mean was **estimated** with the prediction)

The derivations of these equations are found in:

Journel, A.G., 1980. *The lognormal approach to predicting local distributions of selective mining unit grades*. **Mathematical Geology**, 12(4): 285-303.

They are also found in:

Webster, R., and M.A. Oliver. 2008. **Geostatistics for environmental scientists**. 2 ed. John Wiley & Sons Ltd.

Back-transformation of log-SK

For SK, where the mean μ is known:

$$\hat{Z}(\mathbf{x}_0) = \exp \left(\hat{Y}(\mathbf{x}_0) + \frac{\sigma_{SK}^2(\mathbf{x}_0)}{2} \right)$$

The estimate is **increased** always, because the prediction variance σ_{SK}^2 must be positive. This increase is a result of the skewed distribution of the back-transformed prediction variance.

$$\text{var}[\hat{Y}(\mathbf{x}_0)] = \mu^2 e^{\sigma_{SK}^2} \left[1 - \exp \left(-\sigma_{SK}^2(\mathbf{x}_0)/2 \right) \right]$$

The back-transformed prediction variance depends on the mean.

Back-transformation of log-OK

For OK, where the true spatial mean μ is not known:

$$\hat{Z}(\mathbf{x}_0) = \exp \left(\hat{Y}(\mathbf{x}_0) + \frac{\sigma_{\text{OK}}^2(\mathbf{x}_0)}{2} - \psi \right)$$

The estimate is **increased** by the prediction variance σ_{OK}^2 must be positive; but also notice that the **LaGrange multiplier** decreases the prediction.

For OK, it is **not possible to back-transform the variance** in the original units, because μ is unknown (see formula for SK variance back-transformation, above).

Computing confidence limits

It is not necessary to back-transform the prediction variance in order to compute the confidence limits of a prediction in terms of the original variable. Instead:

1. Compute the prediction and its variance in the transformed (e.g., logarithmic) space.
 - This assumes that the random field is Gaussian.
2. Compute the confidence interval in the transformed space.
3. Back-transform the ends of the confidence intervals, simply by exponentiating.

Topic: Kriging with External Drift (KED)

This is a mixed interpolator that includes feature-space predictors that are not geographic coordinates.

The **mathematics are exactly as for UK**, but the *base functions* are different.

In UK, the base functions refer to the **grid coordinates**; these are by definition known at any prediction point.

In KED, the base functions refer to some **feature-space covariate** ...

- ... measured at the sample points (so we can use it to set up the predictive equations) and
- **also known at all prediction points** (so we can use it in the prediction itself).

Base functions for KED

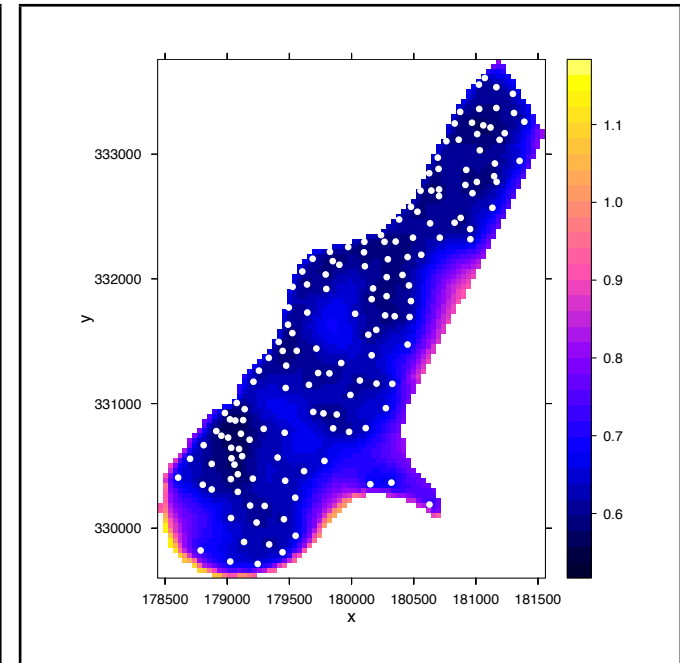
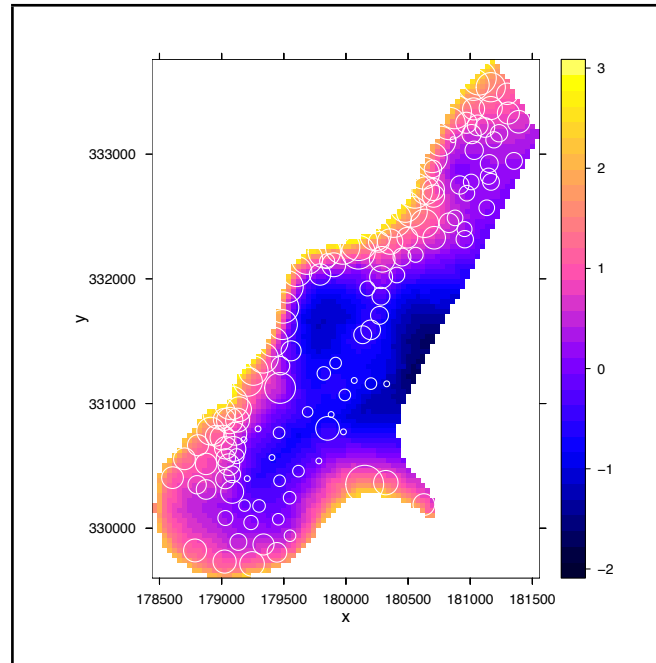
There are two kinds of feature-space covariates:

1. **strata**, i.e. factors. Examples: soil type, flood frequency class
 - Base function: $f_k(\mathbf{x}) = 1$ iff sample or prediction point \mathbf{x} is in class k , otherwise 0 (class indicator variable)
2. **continuous covariates**. Examples: elevation, NDVI
 - Base function: $f_k(\mathbf{x}) = v(\mathbf{x})$, i.e. the value of the predictor at the point.

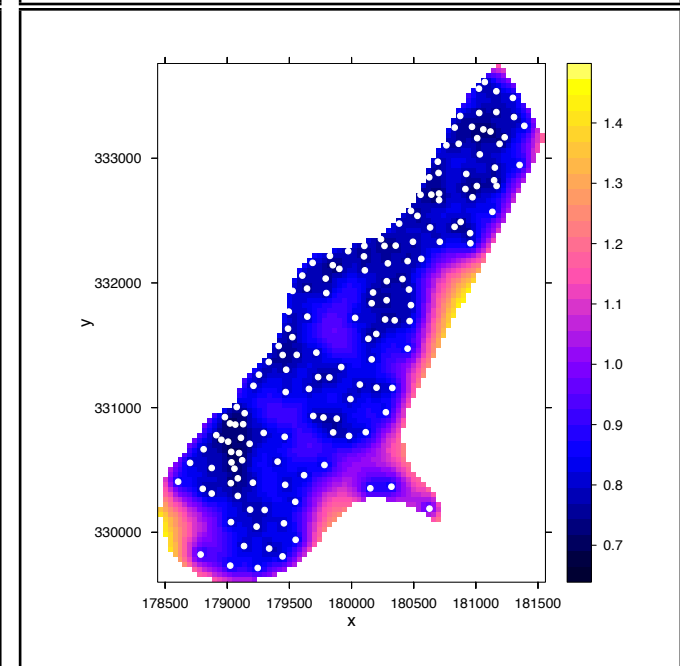
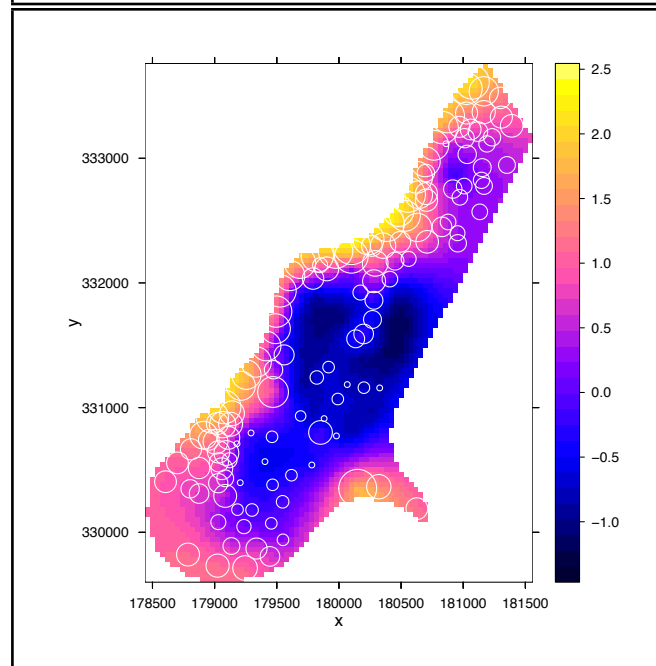
Note that $f_0(\mathbf{x}) = 1$ for all models; this estimates the global mean (as in OK).

Compare KED to OK on the next page; the distance from the river somewhat modifies the predictions. Note the lower kriging variances with KED; this because the residual variogram has lower sill.

KED (distance to river)



OK



Topic: Regression Kriging (RK)

Regression Kriging (RK), also called “kriging after de-trending” models the *trend* (geographic or feature space) and its *residuals* separately.

UK/KED vs. RK (1)

Two methods for incorporating a trend in geographic or feature space:

1. **UK or KED**: Compute trend along with residuals in one UK/KED system
 - UK/KED gives a combined kriging variance
 - Recall: kriging systems are the same, the **base functions** are different
2. **RK**
 - (a) Calculate **trend** (or **strata**) → prediction variance of linear model
 - (b) subtract trend to get **residuals**
 - (c) **model & krige residuals** using Simple Kriging (SK) with known mean of residuals = 0; → kriging variance of residuals
 - (d) **add trend back** to get estimate
 - (e) **sum the two prediction variances** at each point to get the overall error

UK/KED vs. RK (2)

- The coefficients β of the global trend surface are *not* explicit in the UK/KED system; they are *implicit* in the solution and will affect the λ_i at each prediction point
 - * By contrast, in RK we first find the best trend surface, by fitting coefficients β required by the linear model; then we model the residuals with SK
- The UK/KED trend is implicitly computed by **generalized least squares** because the covariance structure is specified
 - * By contrast, the RK trend surface is usually computed without knowing the covariance structure, although that can be estimated by an iterative process.
- The variable being kriged is different with the two approaches: original data values (OK) vs. residuals after fit (RK)
- The λ are different with the two approaches
- There are more ψ (LaGrange multipliers) in UK, and ψ_0 is different in any case
- The RK system does *not* account for the variance of the feature-space model; in the UK/KED system this is built-in.

Naïve RK : Using Ordinary Least Squares (OLS)

- The trend is computed by Ordinary Least Squares (OLS).
- This assumes **no correlation** between residuals (“errors” in statistical terminology):
- Computation:

$$\hat{z}_{\text{RK}}(\mathbf{x}_0) = \sum_{k=0}^p \hat{\beta}_k \cdot q_k(\mathbf{x}_0) + \sum_{i=0}^n \lambda_i z(\mathbf{x}_i)$$

- The two terms are the trend and local interpolation:
 - * p is the degree of the trend surface
 - * the $\hat{\beta}_k$ are estimated by OLS and multiply the covariates q_k evaluated at the prediction point
 - * n is the number of sample points
 - * the λ_i are estimated from the SK system of the residuals and multiply the sampled value of the target variable z at each sample point.

Problems with the naïve RK approach (1)

- Once the trend is removed, the residuals are supposed to be **independent samples** from the **same error distribution**.
- However, we have **evidence** from their variogram that **the regression residuals are in fact correlated** – this is the reason for RK instead of just using the trend.
 - * Goodness-of-fit measure (R^2) of the OLS trend is too optimistic
 - * The trend surface coefficients may be wrong
 - * → the residuals may be wrong
 - * → the variogram of the residuals may be poorly-modelled
 - * → the SK of the residuals may be wrong.
 - * → the prediction may be wrong.
- Four wrongs don't make a right ... what do we do?

Problems with the naïve RK approach (2)

- In practice, if the sample points are **well-distributed** spatially (not clustered), the correlated residuals will be about the same everywhere, so that the OLS fit is satisfactory and naïve RK will give good results.
- *However*, if the sample points are **clustered** in some parts of the map (or feature space), we may well have *mis-estimated the regression coefficients* because we didn't weight the sample points.
 - * N.b. This can happen in non-spatial regression also, where observations in feature space are clustered at certain predictor values.
- In particular, a large number of close-by points with similar values (as is expected by spatial dependence) will “pull” a trend surface or regression towards them
- This will tend also to lower the R^2 (supposed goodness-of-fit) → over-optimistic prediction variance
 - * Thought experiment: add many close-by observations to one of the sample points; if their data values are quite similar (as expected if there is spatial correlation) these will all be well-fit, thereby decreasing the R^2 arbitrarily.

A more sophisticated RK approach: GLS/RK

- Use Generalised Least Squares (GLS), which allows a **covariance structure between residuals** to be **included directly in the least-squares solution** of the regression equation.
- This is a special case of **Weighted Least Squares** (WLS), in this case where samples are weighted according to the spatial structure using the theory of random fields.
- The GLS estimate of the regression coefficients is:

$$\hat{\beta}_{gl\text{s}} = (q^T \cdot C^{-1} \cdot q)^{-1} \cdot q^T \cdot C^{-1} \cdot z$$

where:

- * z is the data vector (observations)
- * q is the *design matrix*;
- * C is the *covariance matrix* of the (spatially-correlated) *residuals*, i.e. deviations from the trend.

GLS/RK (2)

Example: for a first-order trend, q consists of three column vectors:

$$q = [1 \ x \ y]$$

where x is one geographic coördinate and y the other.

Reference: Cressie, N., 1993. *Statistics for spatial data*. John Wiley & Sons, New York.

- The covariance structure must be known (or modelled).
- If there is no spatial dependence among the errors, C reduces to $I\sigma^2$ and the estimate to OLS.

Problem implementing GLS/RK

- The covariance structure refers to the *residuals*, but we can't compute these until we fit the trend ... but we need the covariance structure to fit the trend ... “which came first, the chicken or the egg?”
- In particular, **the covariance structure must be modelled from the variogram of residuals**; but the residuals can only be obtained **after the trend has been computed**.
- These are solved together in UK or KED, but not in RK; here an **iterative** approach is used.

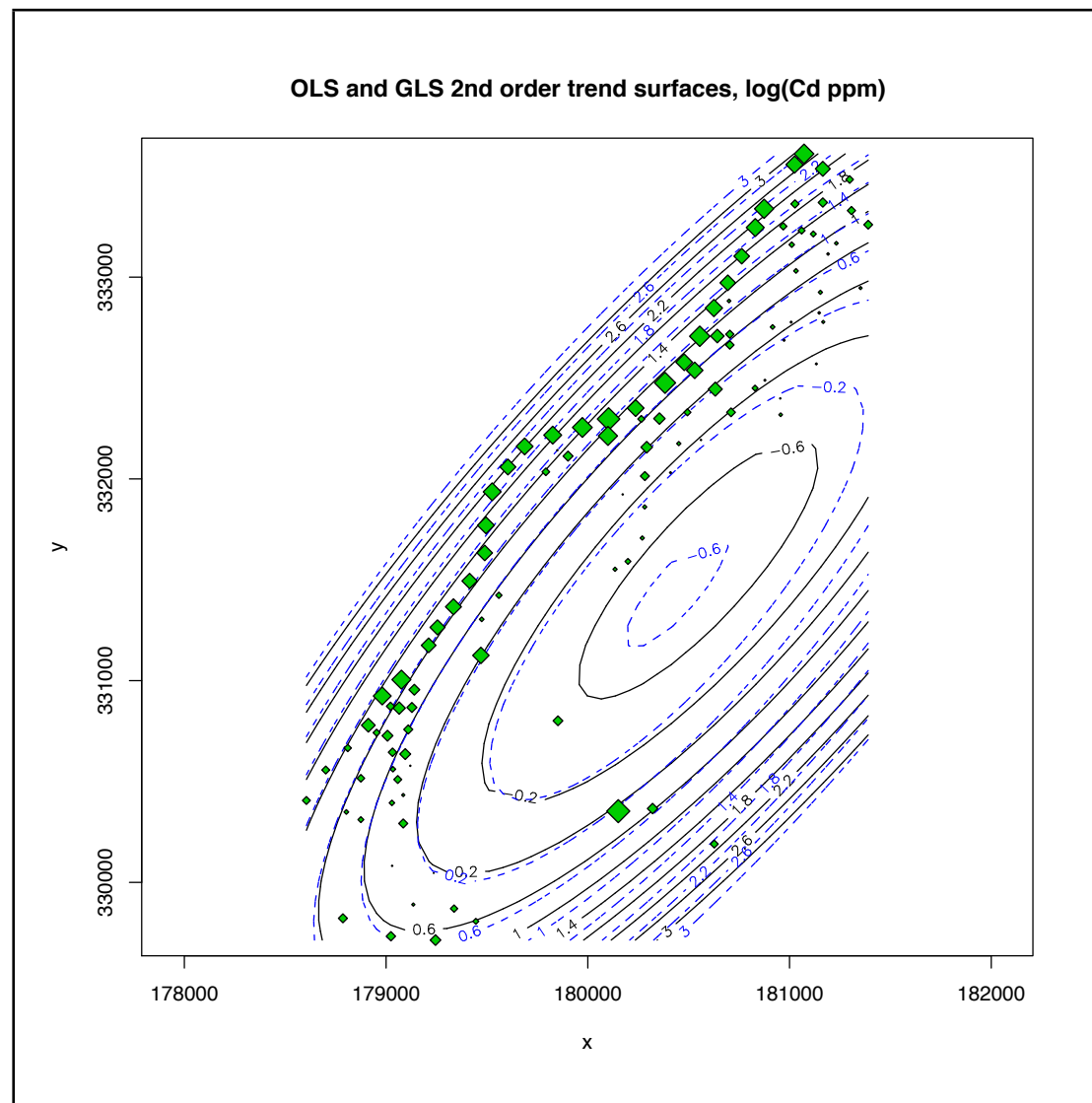
Procedure with GLS/RK

1. Compute the **OLS trend** either on geographic coordinates (as in UK) or on some feature-space covariates (as in KED).
2. Subtract OLS trend from sample points to obtain **OLS residuals**.
3. Model the **covariance function** C of OLS residuals.
4. Compute the **GLS trend** on the geographic- or feature-space covariates using covariance function C to weight the observations.
 - So the trend is modelled **twice**, first with OLS just to get the residuals to model the spatial covariance, and second with GLS using this modelled covariance.
5. Subtract GLS trend from sample points to obtain **GLS residuals**.
6. Model the **semi-variance function** $\gamma(\mathbf{h})$ of GLS residuals ...

7. **Predict GLS residuals** at prediction points by SK ($\mu \equiv 0$)
8. **Predict GLS trend** at prediction points from the GLS regression equation
9. **Add** predicted GLS trend and predicted GLS residuals at prediction points
 - → **final prediction**
10. **Add** GLS prediction variance and SK kriging variance at prediction points
 - → **final kriging variance**

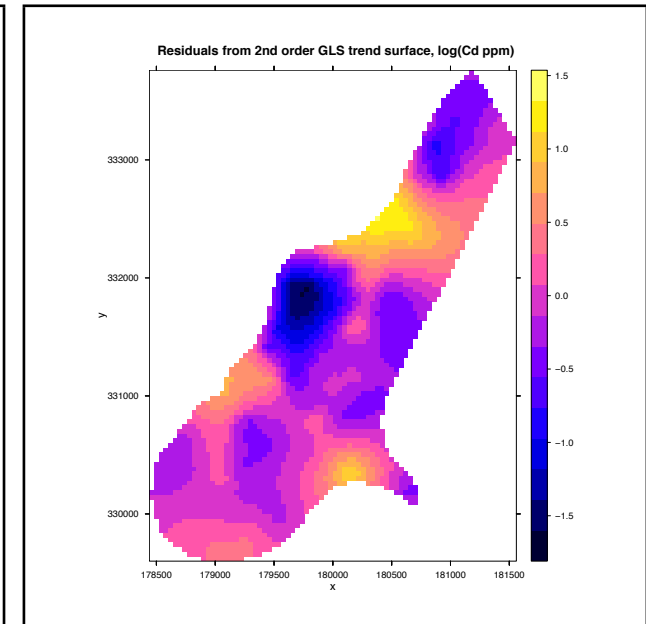
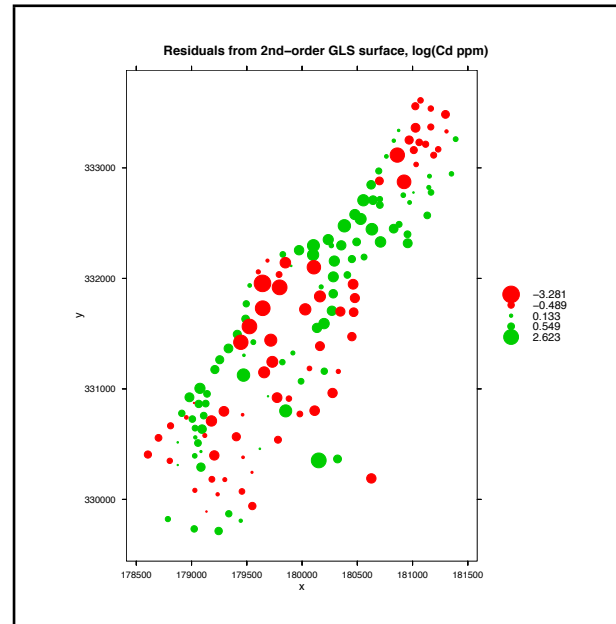
Visualizing GLS/RK (1)

OLS and GLS trend surfaces compared

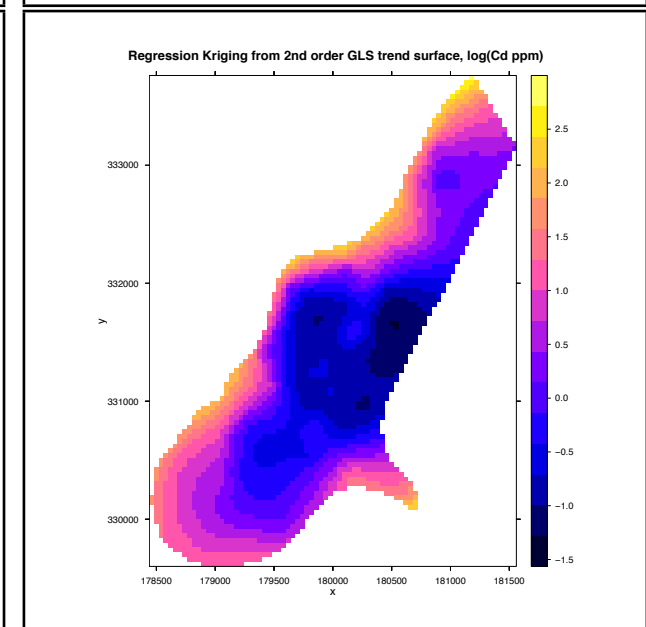
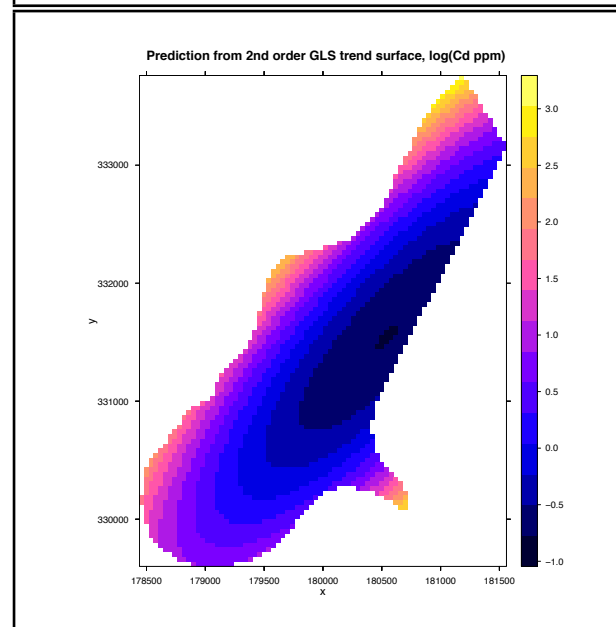


Visualizing GLS/RK (2)

Residuals, SK

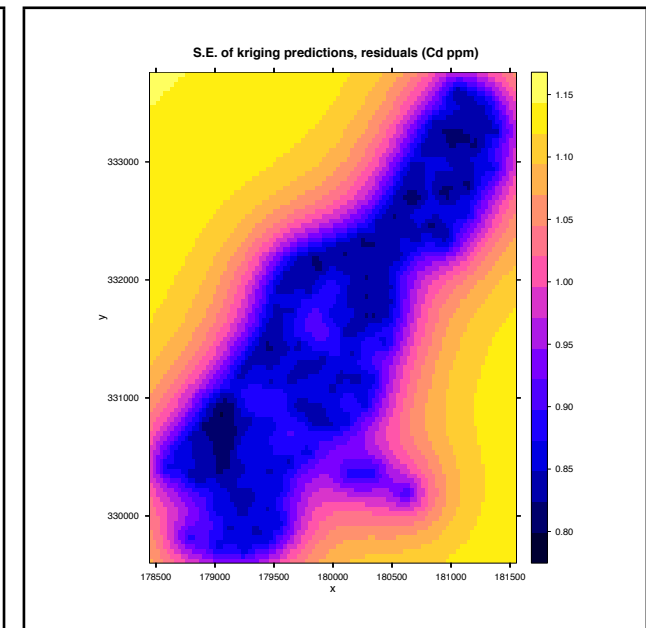
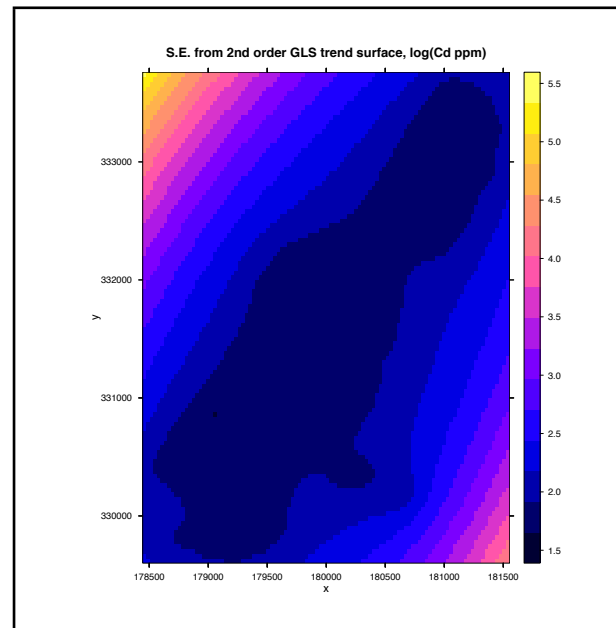


TS, RK

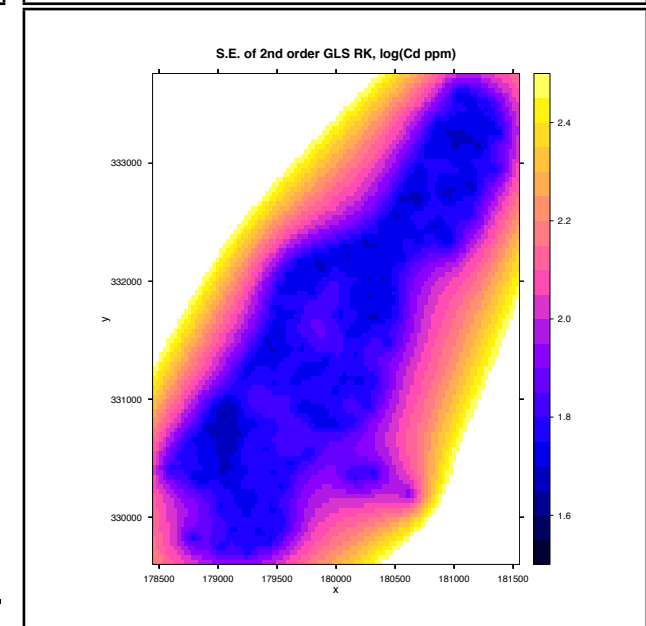


Visualizing GLS/RK (3)

Trend surface, SK S.E.



Regression Kriging S.E.



An even more sophisticated approach – REML

Generalized Least Squares (GLS) is an iterative procedure (“chicken-and-egg”), so there is no assurance we know the correct covariance structure among the regression residuals.

It does provide unbiased estimates of the fixed effects (i.e., regression coefficients) but the variance parameters of these (e.g., standard errors of coefficients, goodness-of-fit) are biased. For large datasets this bias becomes small, so the above approach in practice usually performs well.

A sophisticated approach is to use Restricted Maximum Likelihood (REML) to estimate both the spatial structure and regression parameters together; see the clear explanation in:

Lark, R. M., & Cullis, B. R. (2004). *Model based analysis using REML for inference from systematically sampled data on soil*. European Journal of Soil Science, 55(4), 799-813.

Topic: Stratified Kriging (StK)

Stratified Kriging (StK) first stratifies an area according to some classified factor, and then predicts within each stratum separately.

The per-stratum results are then combined for a final map.

StK vs. KED with a classified predictor

One variant of **Kriging with External Drift** is for the “drift” to be represented by some classifying factor which is mapped as polygons.

Examples: soil types; lithological units; political subdivisions

In this form of KED, the **mean** is expected to differ among strata.

This restores **first-order** stationarity to the **residuals**, which are then modelled and predicted by Simple Kriging.

StK also allows the **spatial structure** within strata to differ.

That is, it does not assume **second-order** stationarity of the residuals.

Example

On a steep hillslope the target variable 'depth to a root-restricting layer' may be not only ...

- **less** (shallower soils) → different **expected value**, but also
- **variable at shorter ranges** → reduced **range** of the variogram model, and perhaps
- **more variable overall** → higher **sill** of the variogram model

... than the depth on a gently-sloping area.

Even the **variogram form** may differ between strata.

Example application: Stein, A., Hoogerwerf, M., & Bouma, J., 1988. *Use of soil-map delineations to improve (co)kriging of point data*. **Geoderma**, 43, 163–177.

Procedure for StK

1. Delineate contrasting areas as strata (i.e., a polygon map);
2. Model within **each stratum separately**;
 - (a) Empirical variogram
 - (b) Fitted variogram model
3. Use these models to map by OK **within each stratum** separately;
4. **Combine** the per-stratum maps into a single map.

Applicability of StK

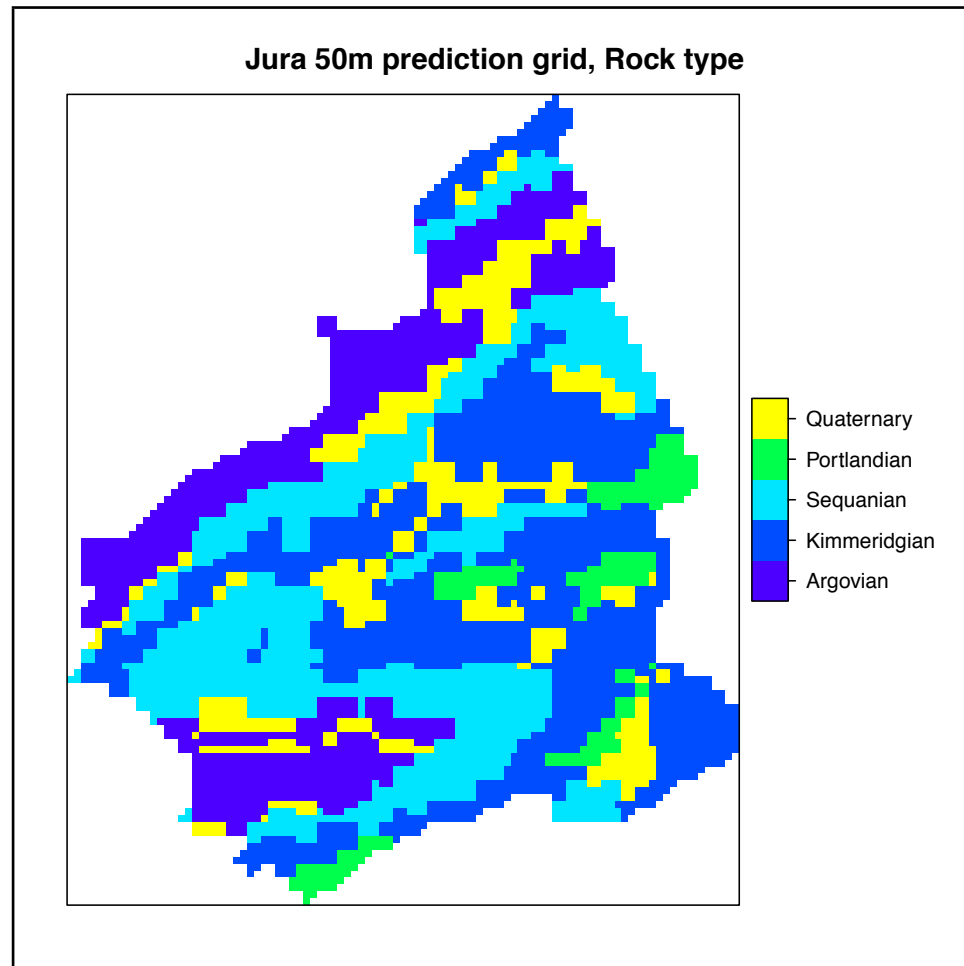
This approach is not much used, because of the following difficulties:

- The requirement to model per-stratum variograms
 - * the number of sample points is often quite small when the area is stratified;
 - * polygons of the stratum may be small and widely-spaced, so the variogram is limited to short-range and may not reach a sill or have a clear model form within that range;
- Predictions use only the known points within each stratum, so are based on limited information;
- If the polygons are widely-spaced, only points within a polygon have large weight, so the prediction is strongly local;
- **Abrupt changes** in both values and prediction variances at stratum boundaries.

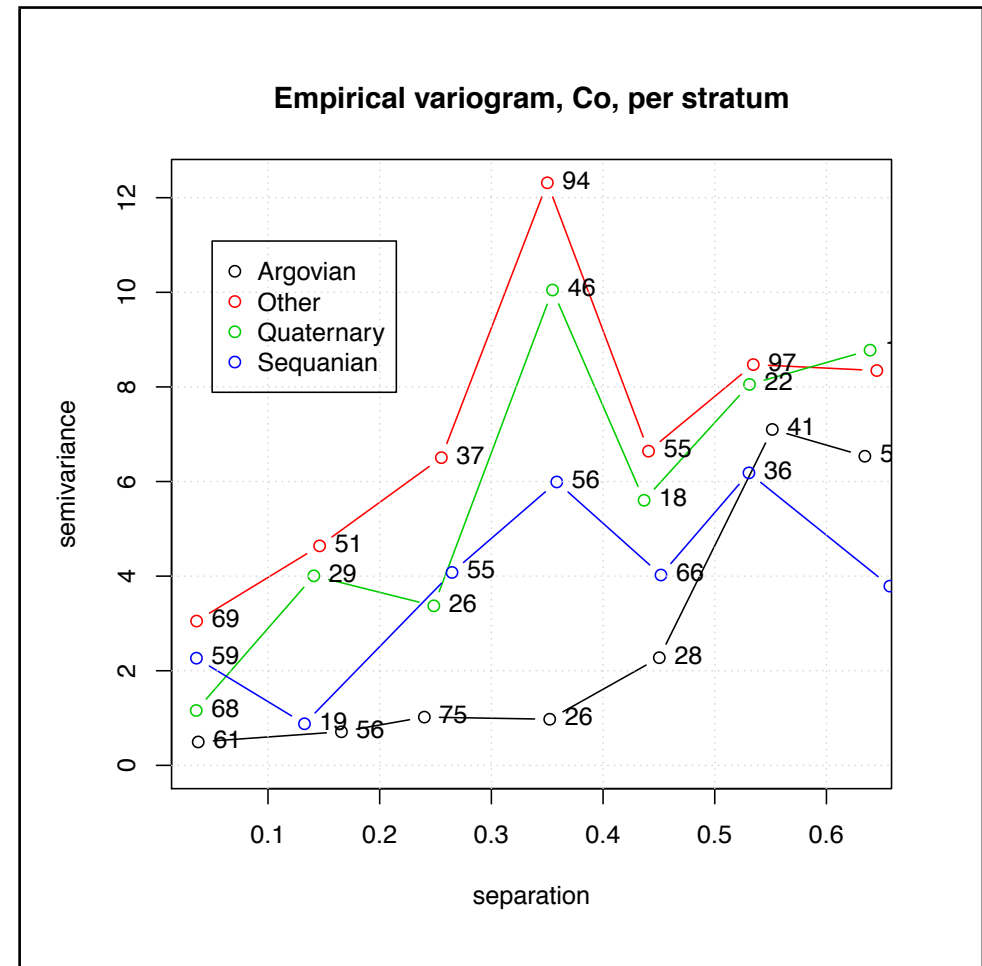
Note this last may be desirable if the strata are strongly-contrasting. Example: depth to bedrock at the interface of a steep hillside and an alluvial plain.

Example: Jura Co concentration, rock type strata

Jura rock types

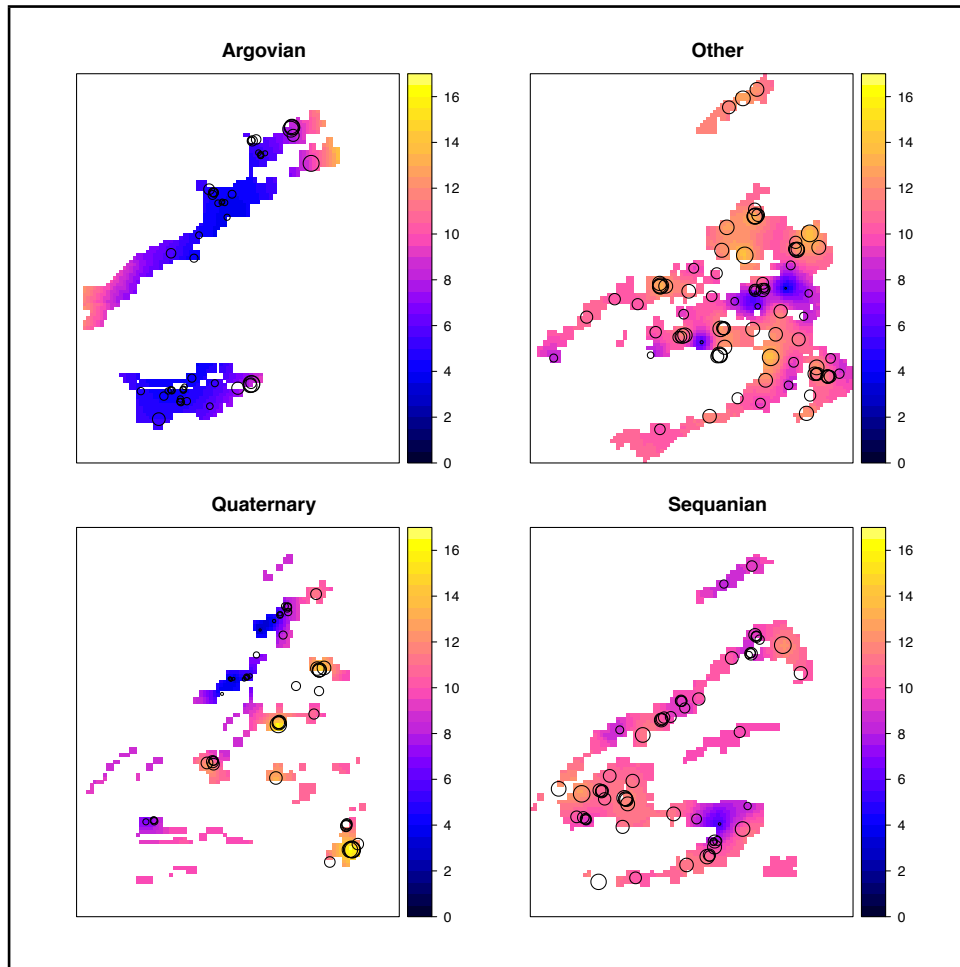


Per-stratum empirical variograms

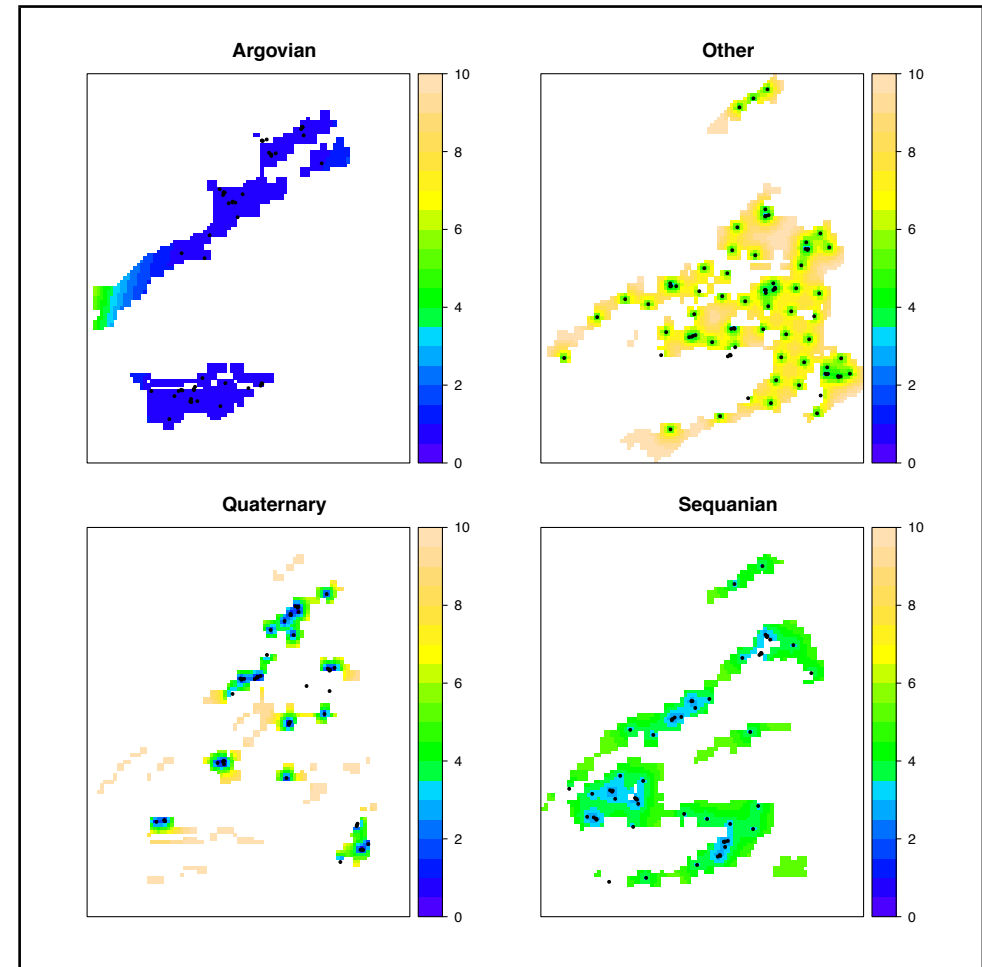


Stratified Kriging results – per-stratum

Per-stratum StK predictions

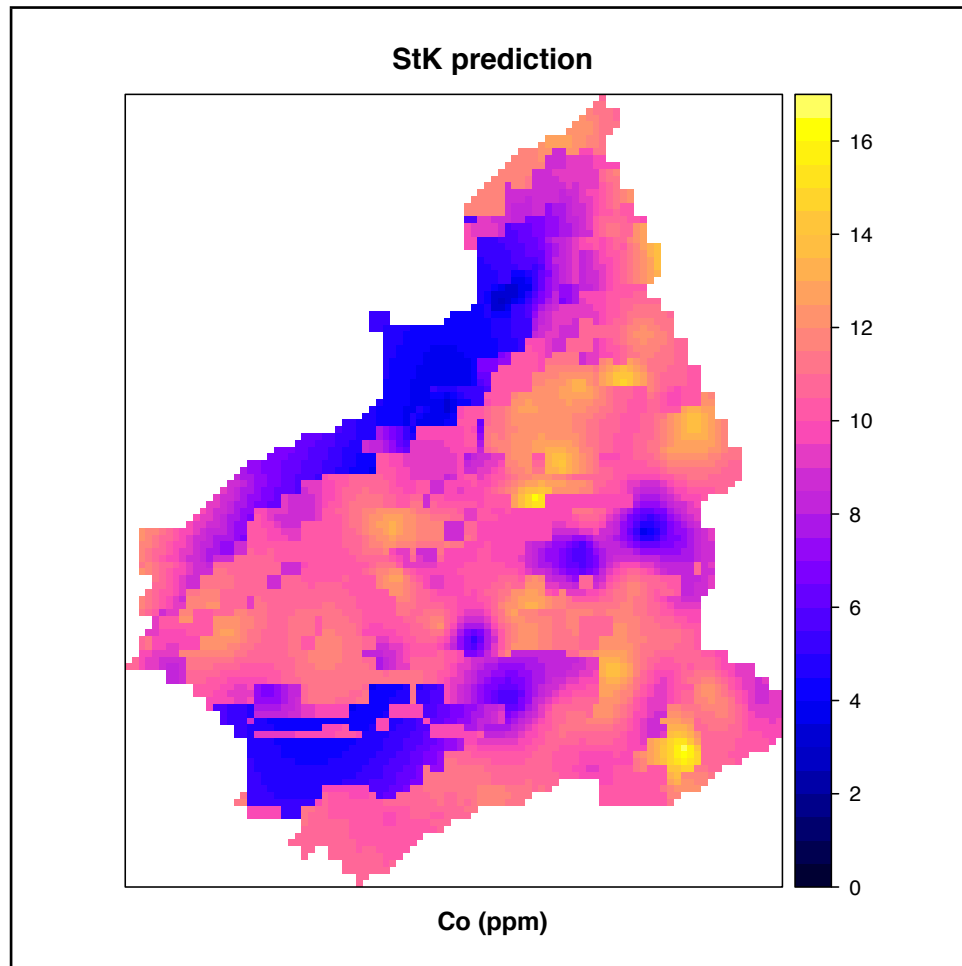


Per-stratum StK prediction variances

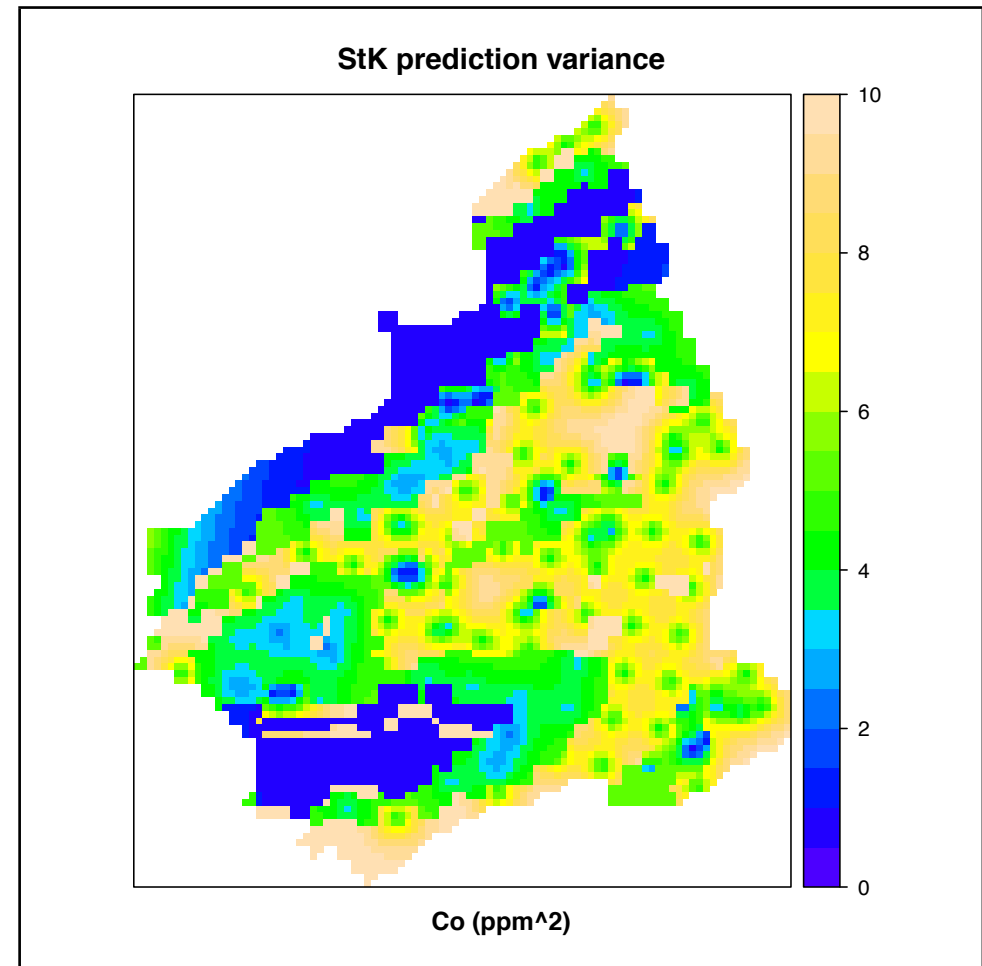


Stratified Kriging results – combined

Combined StK prediction

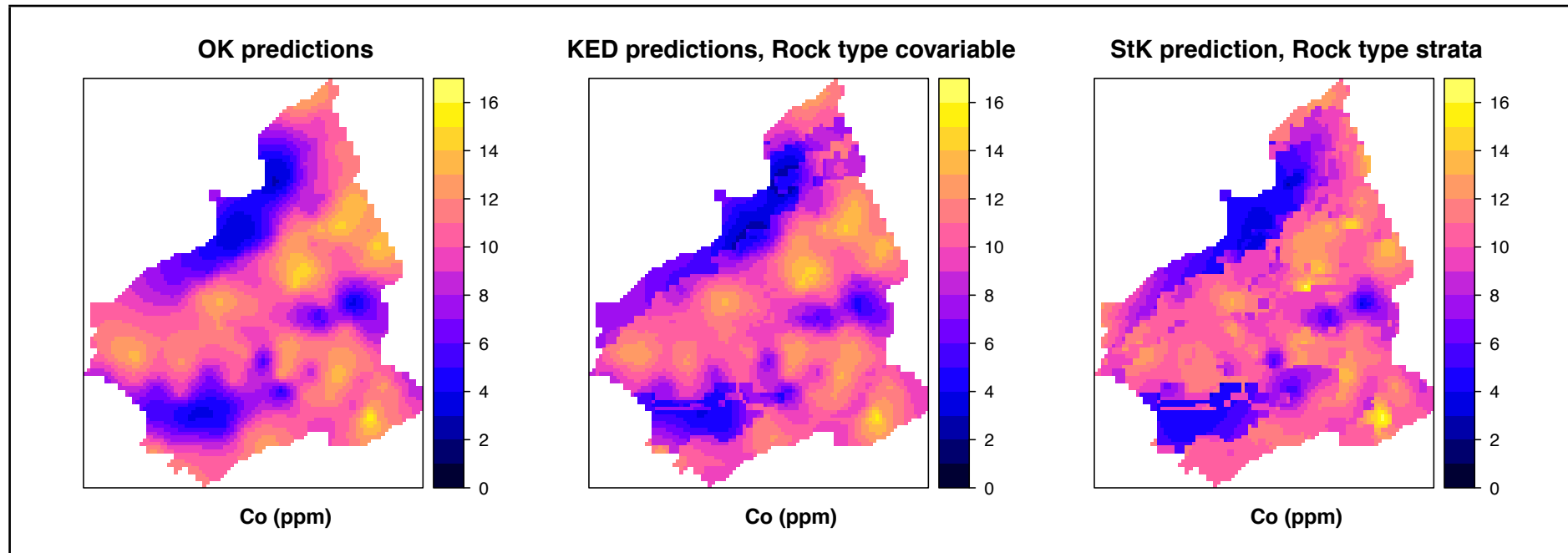


Combined StK prediction variances



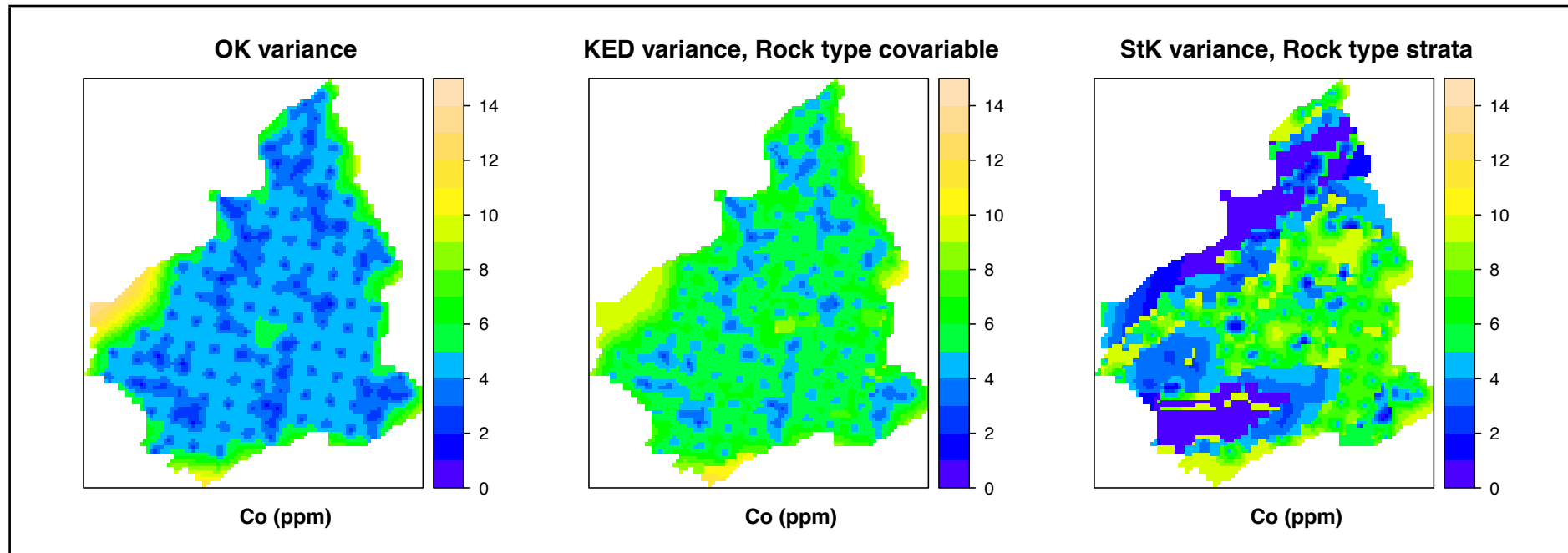
Comparing OK, KED, StK (1)

Comparing OK, KED, StK predictions



Comparing OK, KED, StK (2)

Comparing OK, KED, StK prediction variances



Topic: Cokriging (CoK)

Often we have several **related sets** of point observations of **different attributes**:

- **co-located**: all attributes at all locations
 - * e.g.: geochemical soil or rock samples with many measured elements
- **partially** co-located: some same locations and all attributes; but additional locations with only some attributes
 - * e.g.: soil moisture at instrumented locations, clay and organic matter at many others
 - * the **under-sampled** attribute is usually the **target**
- **disjunct** locations: each attribute is collected at different location sets

Coregionalization and Cokriging

These two concepts are related in the same way as the *theory of regionalized variables* and (univariate) *kriging*:

1. **Coregionalization** is a theoretical model of how several variables **spatially co-vary**; this is used for ...
2. **Cokriging** (CK), which is a method of using:
 - supplementary information on a **co-variable** ...
 - ... to improve the **prediction** of a **target variable** ...
 - or to **predict several related variables** at the same time

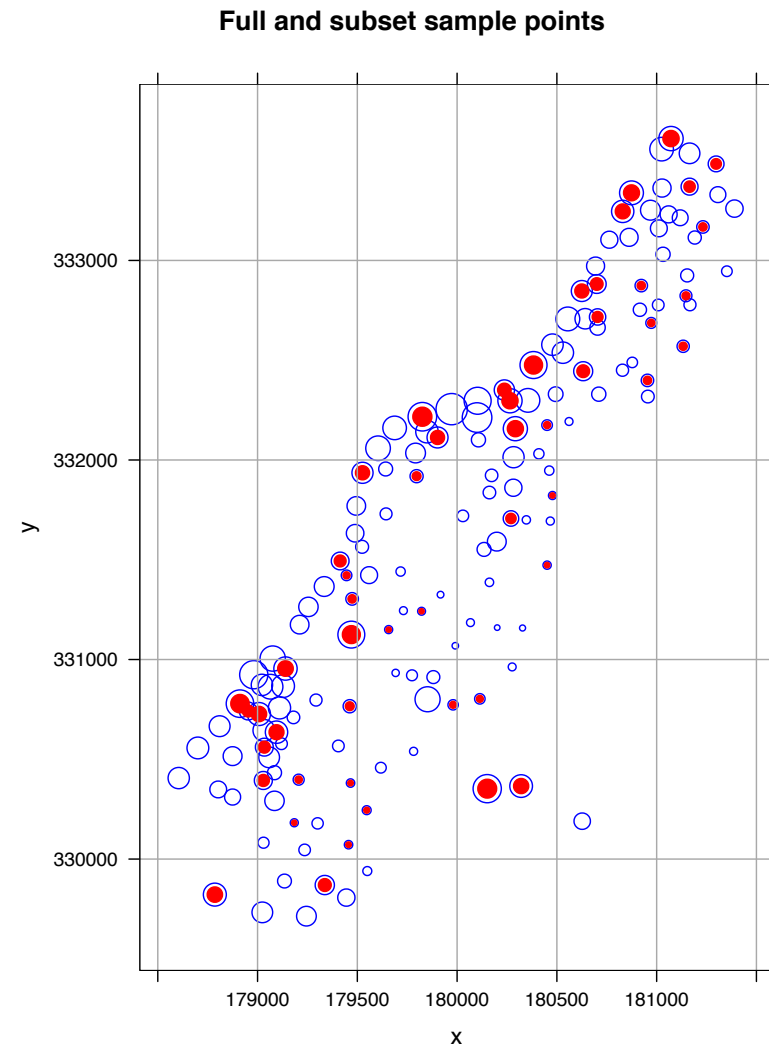
Steps in Cokriging

1. Extend the **theory of regionalised variables** to several variables
2. Examine **spatial structure** with direct and cross-**variograms**
 - show that there is a **multivariate spatial cross-correlation** ...
 - ... as well as the **univariate spatial auto-correlation**
 - These are then called **co-regionalised** variables
3. **Model** this spatial structure (direct and cross-variogram models)
4. Build the **cokriging system** of equations
5. Use it to **predict** at unsampled locations.

When to consider CoK?

- When the target attribute is undersampled, and:
 - * there are one or more other attributes **co-variables** measured at more locations (maybe some co-located)
 - * there is a strong **spatial cross-correlation** between the target and co-variables (see below)
 - * generally the under-sampled target is **more expensive**/difficult to measure
 - e.g.: a set of lab. determinations of soil organic matter vs. a larger set of spectroscopic measurements
- When several targets are all to be mapped, and:
 - * they are closely-related in feature space (**strong feature-space correlation**)
 - * there is a strong **spatial cross-correlation** between the target and co-variables (see below)
 - e.g.: a set of closely-related geochemical elements in soil/rock samples

Example of under-sampling



Two sampling intensities, **some co-located observations**

Main difficulty of CoK

- Building an authorized model of spatial **co**variance that results in a **positive-definite CoK system**
- This because the covariables may have a different spatial structure
 - * they may each be auto-correlated, but may have different structures
 - * so the **linear model of co-regionalization** can not be used
 - * other models are quite difficult to formulate

Alternatives to CoK with covariables

If the covariable(s) is (are) known at **all** prediction locations, use Kriging with External Drift (**KED**) or Regression Kriging (**RK**)

These **do not require a model of spatial covariance**, they just use feature-space correlation.

E.g.: soil clay content (measured at points) is the target; MODIS day-night temperature differences for some days after a rain are the covariable(s); MODIS is available over the whole prediction grid.

This is a typical example when **remote-sensing** data are covariables.

Reference: Zhao, M.-S., D.G. Rossiter, D.-C. Li, Y.-G. Zhao, F. Liu, and G.-L. Zhang. 2014. *Mapping soil organic matter in low-relief areas based on land surface diurnal temperature difference and a vegetation index*. **Ecological Indicators** 39: 120–133.

Correlation in feature space vs. geographic space

- In **feature space**, the correlation is between the two variables at **co-located points** only
- In **geographic space**, the correlation is between the two variables at **point-pairs separated by a distance** and summarized in **distance classes**
 - * The absolute value of the correlation is expected to be greatest if points are co-located, and then decrease with separation – **if the variables are co-regionalized**
 - * Otherwise there should be no relation between variables that are not co-located

Variograms

Two types of variograms:

1. **direct**: single regionalised variables, one variogram *per variable*
 - compute in the usual way for a univariate empirical variogram
2. **cross**: *per pair* of regionalised variables

Empirical cross-variogram

This is analogous to the standard Matheron estimator for the direct variogram.

For the two variables u and v with values z_u and z_v respectively:

$$\hat{\gamma}_{uv,\mathbf{h}} = \frac{1}{2m(\mathbf{h})} \sum_{i=1}^{m(\mathbf{h})} \{z_u(\mathbf{x}_i) - z_u(\mathbf{x}_i + \mathbf{h})\} \{z_v(\mathbf{x}_i) - z_v(\mathbf{x}_i + \mathbf{h})\}$$

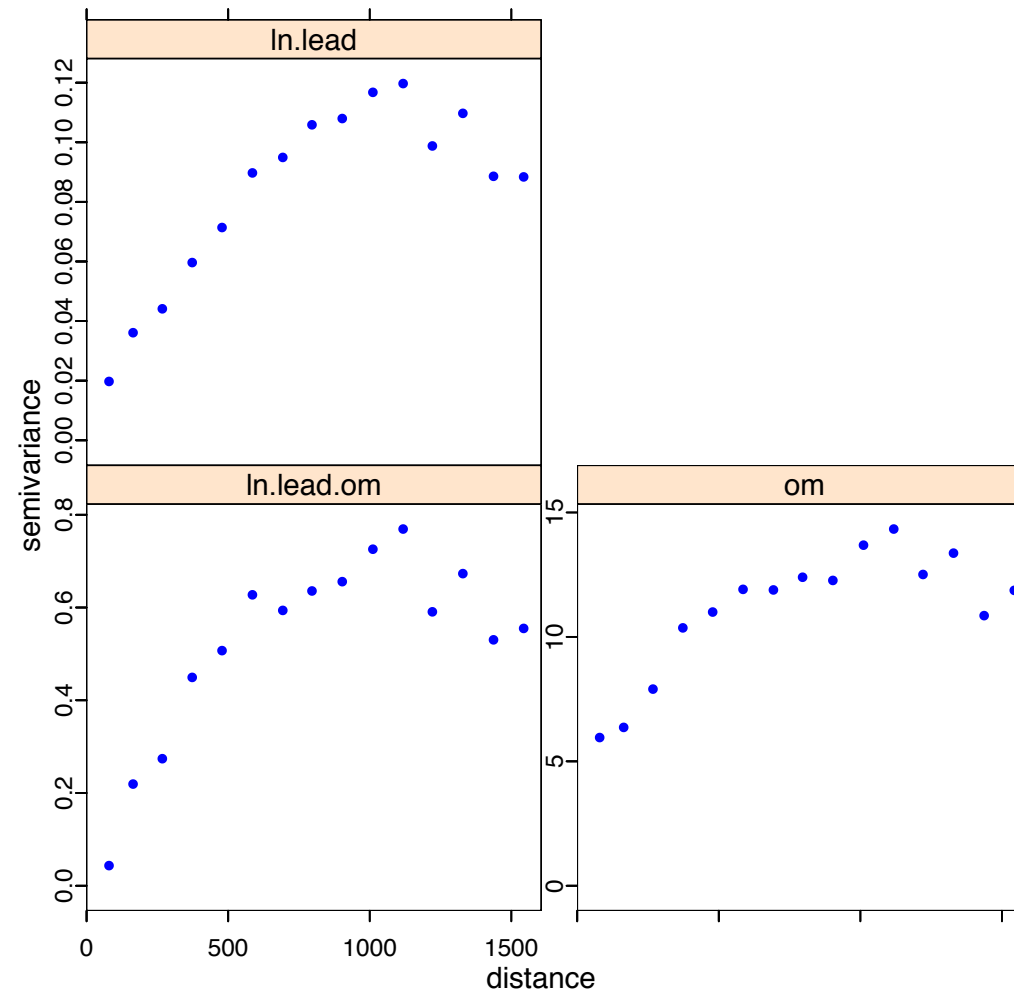
where $m(\mathbf{h})$ is the number of point-pairs separated by vector \mathbf{h} ; for an omnidirectional variogram this is the distance class

In words: if high differences between point-pairs of one variable are positively associated with high differences between point-pairs of the other variable, they will have a *high positive cross-correlation*.

This can also be a **negative** cross-correlation!

If the differences are randomly associated, there will be on average no spatial cross-correlation in this distance class.

Empirical direct- and cross-variograms



Here a strong **positive cross-correlation** (soil Pb vs. organic matter)

Modelling the variograms

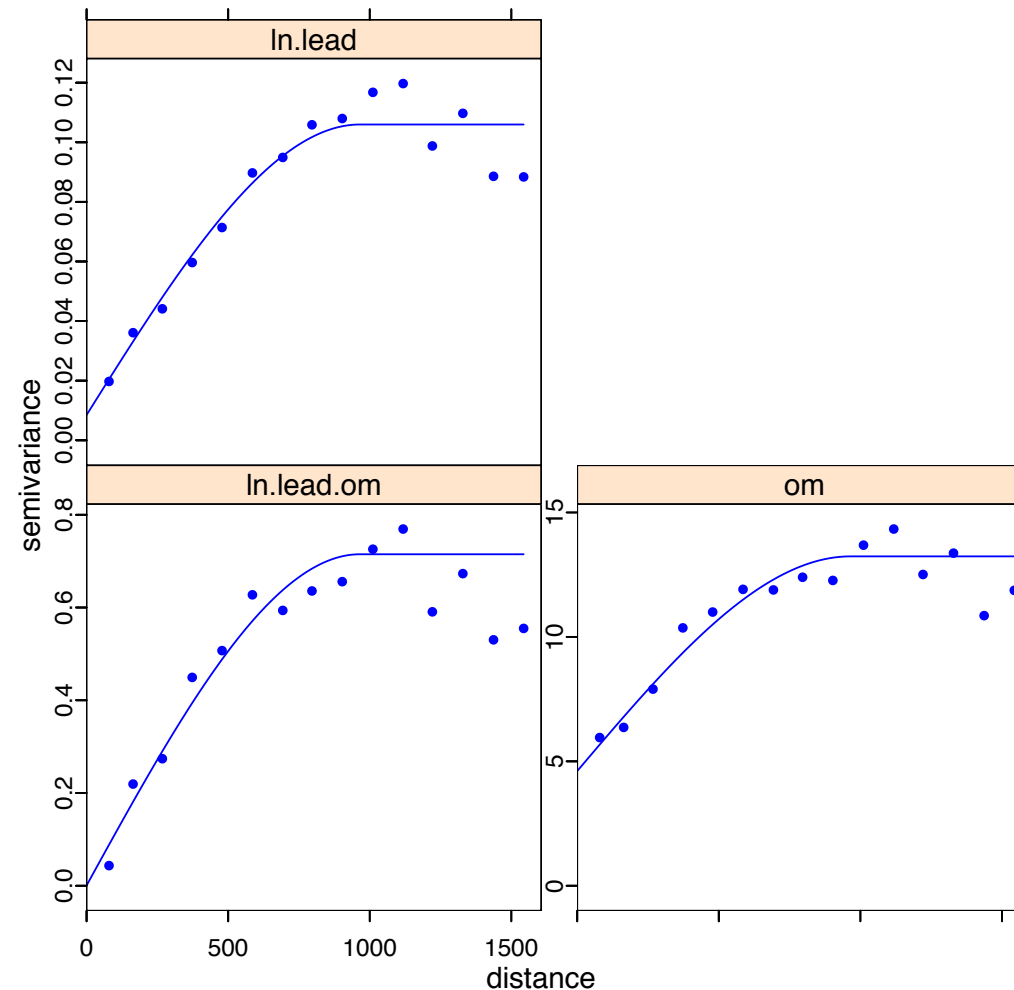
The **direct and cross-variograms must be modelled together**, with some restrictions to ensure that the resulting CK system can be solved.

The simplest way to ensure this is to assume a **linear model of co-regionalisation**: *all variograms are linearly related*

- Same **model**, same **range** (i.e., same spatial structure)
- May have different **sills** (i.e., overall variability)
- May have different **nuggets** (i.e., uncertainty at sample point)

Other models are possible but much more complicated to estimate, while ensuring a **positive definite** cokriging system

Direct- and cross-variograms with fitted models



Reasonable fit, although range of the cross-variogram seems shorter. Fits well at shorter ranges, where most of kriging weights are determined.

Co-Kriging : Prerequisites

1. Two point data sets, usually with some observations co-located:
 - (a) the **target variable** z at locations $\mathbf{x}_1 \dots \mathbf{x}_{N_z}$
 - (b) the **co-variable** w at locations $\mathbf{y}_1 \dots \mathbf{y}_{N_w}$
2. **Spatial structure in both variables separately** (i.e. non-nugget **variograms**);
3. **Spatial structure between the variables (cross-variograms)**; this can be either a positive or negative spatial correlation;
4. Certain restrictions on the joint spatial structure.

(N.b. can use more than one co-variable but we will only cover the case where there is only one.)

The Co-kriging predictor

Prediction of the target variable at an unknown point \mathbf{x}_0 is computed as the **sum of two weighted averages**:

1. one of the N_z sample values of the **target variable** z ; and
2. one of the N_w sample values of the **co-variable** w

$$z(\mathbf{x}_0) = \sum_{i=1}^{N_z} \lambda_i z(\mathbf{x}_i) + \sum_{j=1}^{N_w} \mu_j w(\mathbf{y}_j)$$

We want to find the weights λ (for the target variable) and μ (for the co-variable) which minimize the prediction variance; this will then be the BLUP.

Unbiasedness

Two conditions:

$$\sum_{i=1}^{N_z} \lambda_i = 1$$

$$\sum_{j=1}^{N_w} \mu_j = 0$$

That is, the weights from the target variable sum to 1, as in OK.

The weights from the co-variable sum to 0, so that there is no overall effect (n.b., this variable may well have different units).

The Co-kriging variance

Suppose there are V variables (target and co-variables), which are indexed from $1 \dots l$; each has n_l observations.

Variable u is the target, one of the V .

we want to minimize the variance.

$$\sigma_u^2 = \sum_{l=1}^V \sum_{j=1}^{n_l} \lambda_{jl} \gamma_{ul}(\mathbf{x}_j, \mathbf{x}_0) + \psi_u$$

So the semivariance of all observations of all variables with the point to be estimated is **minimized**.

There is one of these equations for each variable.

The Co-Kriging system (1)

Solve: $\mathbf{A}_C \lambda_C = \mathbf{b}_C$, where \mathbf{A}_C is built up from the direct and cross-semivariances of the sample points (both target and co-variable):

$$\mathbf{A}_C = \begin{bmatrix} \Gamma_{zz} & \Gamma_{zw} & \mathbf{1} & \mathbf{0} \\ \Gamma_{wz} & \Gamma_{ww} & \mathbf{0} & \mathbf{1} \\ \mathbf{1}^T & \mathbf{0}^T & 0 & 0 \\ \mathbf{0}^T & \mathbf{1}^T & 0 & 0 \end{bmatrix}$$

The vectors of 1's and 0's control which of the semi-variances are included in each equation.

The Γ are the matrices of semi-variances (next slide)

By linear algebra, the solution is: $\lambda_C = \mathbf{A}_C^{-1} \mathbf{b}_C$

The Co-Kriging system (2)

The Γ are the matrices of semi-variances:

1. Γ_{zz} (dimension $N_z \times N_z$) between the locations of the **target** variable, computed from the **direct** variogram for z ;
2. Γ_{ww} (dimension $N_w \times N_w$) between the locations of the **co**-variable, computed from the **direct** variogram for w ;
3. Γ_{wz} (dimension $N_w \times N_z$) between the locations of the **target** and **co**-variables, computed from the **cross**-variogram;
4. $\Gamma_{zw} = \Gamma_{wz}^T$

Example of a cross-variable matrix

$$\Gamma_{zw} = \begin{bmatrix} \gamma_{zw}(\mathbf{x}_1, \mathbf{y}_1) & \gamma_{zw}(\mathbf{x}_1, \mathbf{y}_2) & \cdots & \gamma_{zw}(\mathbf{x}_1, \mathbf{y}_{N_w}) \\ \gamma_{zw}(\mathbf{x}_2, \mathbf{y}_1) & \gamma_{zw}(\mathbf{x}_2, \mathbf{y}_2) & \cdots & \gamma_{zw}(\mathbf{x}_2, \mathbf{y}_{N_w}) \\ \vdots & \vdots & \cdots & \vdots \\ \gamma_{zw}(\mathbf{x}_{N_z}, \mathbf{y}_1) & \gamma_{zw}(\mathbf{x}_{N_z}, \mathbf{y}_2) & \cdots & \gamma_{zw}(\mathbf{x}_{N_z}, \mathbf{y}_{N_w}) \end{bmatrix}$$

The Co-Kriging system (3)

$$\lambda_{\mathbf{C}} = \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_{N_z} \\ \mu_1 \\ \vdots \\ \mu_{N_w} \\ \psi_z \\ \psi_w \end{bmatrix} \quad \mathbf{b}_{\mathbf{C}} = \begin{bmatrix} \gamma_{zz}(\mathbf{x}_1, \mathbf{x}_0) \\ \vdots \\ \gamma_{zz}(\mathbf{x}_{N_z}, \mathbf{x}_0) \\ \gamma_{wz}(\mathbf{y}_1, \mathbf{x}_0) \\ \vdots \\ \gamma_{wz}(\mathbf{y}_{N_w}, \mathbf{x}_0) \\ 1 \\ 0 \end{bmatrix}$$

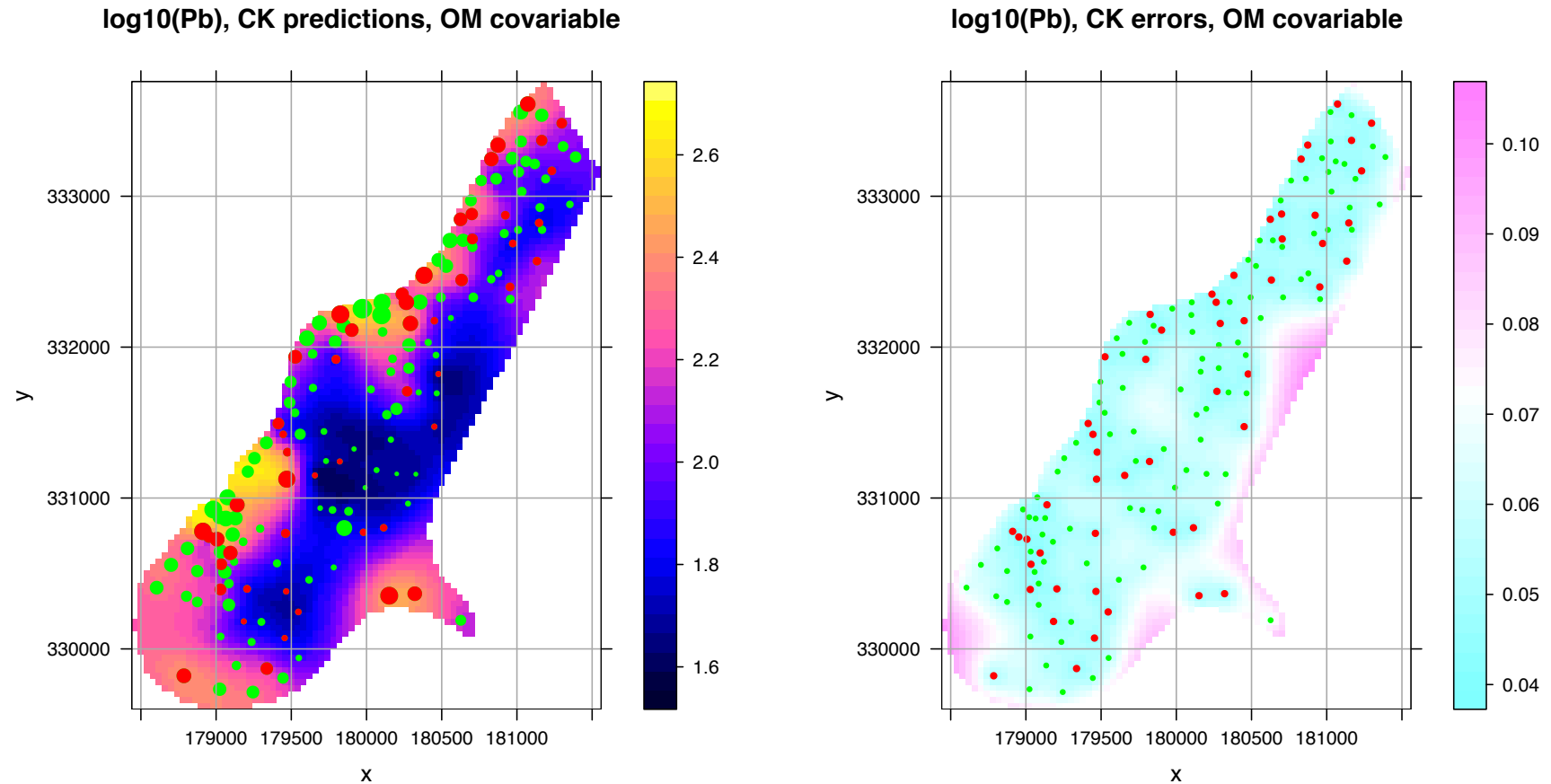
Note that $\mathbf{b}_{\mathbf{C}}$ uses **direct** semivariances between the prediction point and the **target** variable, but **cross**-semivariances between the prediction point and the **co**-variable. Note also one LaGrange multiplier for each variable.

The **cokriging prediction variance** is then:

$$\hat{\sigma}_z^2(\mathbf{x}_0) = \mathbf{b}_{\mathbf{C}}^T \lambda_{\mathbf{C}}$$

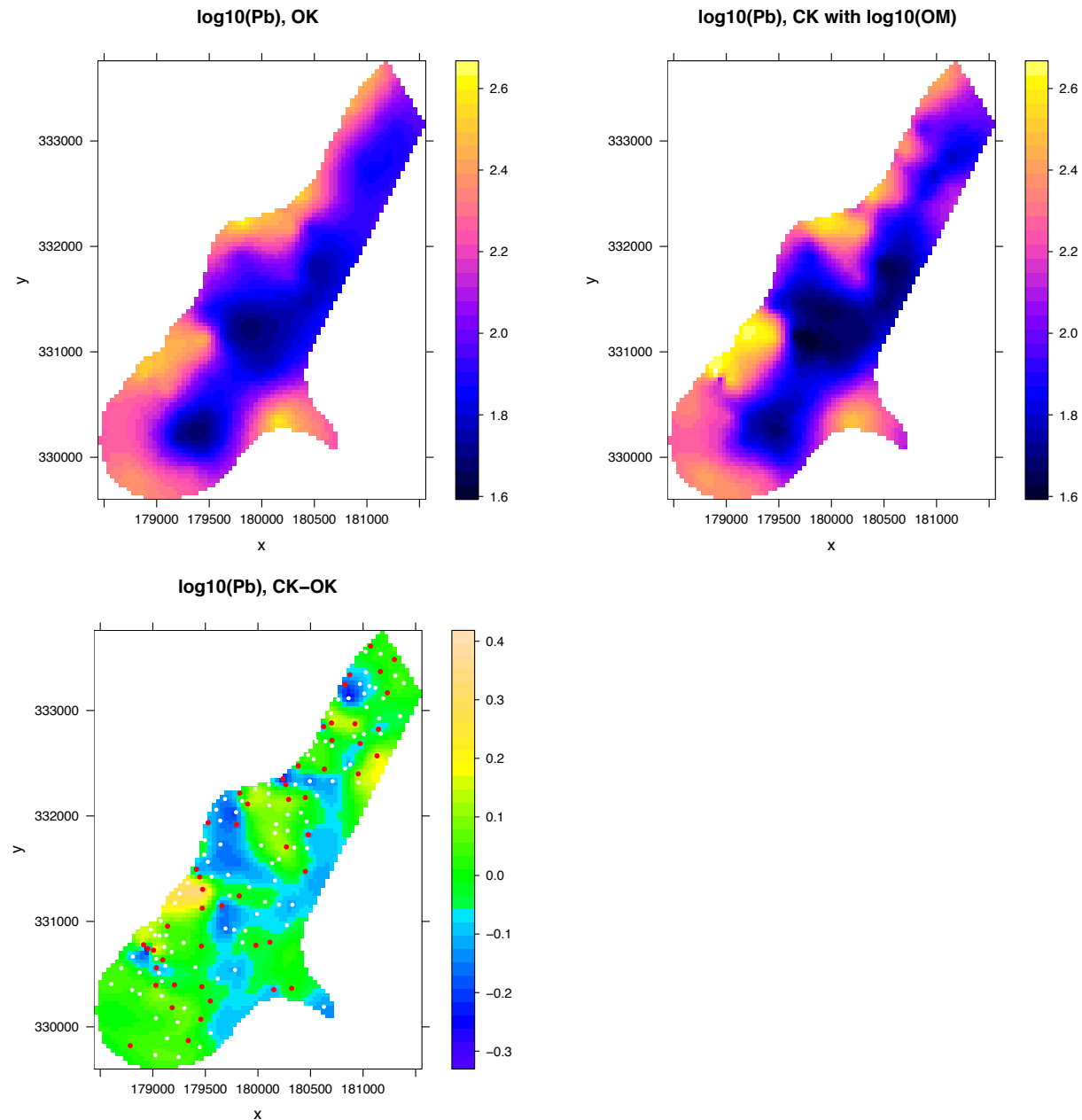
Example cokriging results

Co-kriging predictions and prediction variance, $\log_{10}(\text{Pb})$



(subsample points red, extra points green)

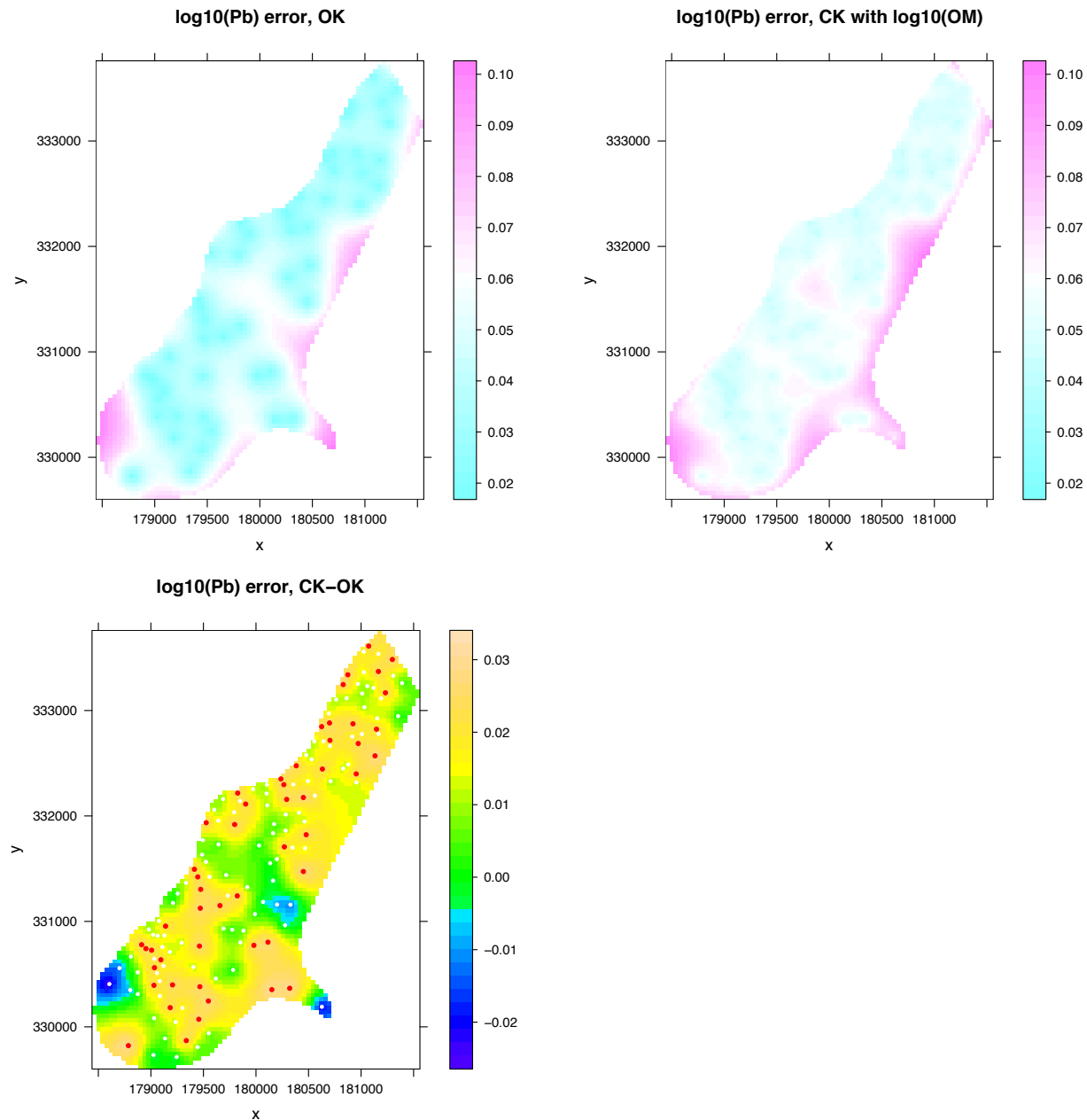
Difference between OK and CK predictions



To check your understanding ...

Q10 : *What are the differences between the **kriging predictions** of OK and CoK? Explain the spatial pattern of these.* *Jump to A10 •*

Difference between OK and CK prediction variances



To check your understanding ...

Q11 : *What are the differences between the **kriging variances** of OK and CoK? Explain the spatial pattern of these.* *Jump to A11 •*

References for Cokriging

The mathematics of this topic is necessarily quite involved, and there are differences in notation, theory, and practice.

Goovaerts, P., 1997. *Geostatistics for natural resources evaluation*. Oxford University Press, Oxford and New York; Chapter 6

Webster, R., & Oliver, M. A., 2001. *Geostatistics for environmental scientists*. Wiley & Sons, Chichester; Chapter 9

Papritz, A., & Stein, A., 1999. *Spatial prediction by linear kriging*. In: A. Stein, F. v. d. Meer & B. G. F. Gorte (Eds.) *Spatial statistics for remote sensing* (pp. 83-113). Dordrecht: Kluwer Academic.

Isaaks, E.H. and Srivastava, R.M., 1990. *An introduction to applied geostatistics*. Oxford University Press, New York.

Rossiter, D.G., 2012. *Technical Note: Co-kriging with the gstat package of the R environment for statistical computing*. ITC, Enschede.

http://www.itc.nl/personal/rossiter/teach/R/R_ck.pdf

Answers

Q1 : *Consider a regular grid over some area; measure the data values and compute the ordinary average.*

If there is spatial dependence, what would happen to this average if we now make many observations near one of the grid points (i.e. a cluster), and not the others? •

A1 : *The average will get closer and closer to the average of the points in the cluster, which, because of spatial autocorrelation, will not necessarily be the average over the whole area. Return to Q1* •

Answers

Q2 : *How is the covariance matrix \mathbf{C} computed?* •

A2 : *This is derived from the modelled **covariance function**: each matrix entry c_{ij} is computed from the **separation** h_{ij} between the two observations z_i and z_j (in the simple case, the distance between them), which is then used as an argument to the covariance function.*

For example, for an exponential model: $c_{ij} = ce^{(-h_{ij}/a)}$

Return to Q2 •

Q3 : *How is the covariance vector \mathbf{c}_0 computed?* •

A3 : *As for the covariance matrix, but the separations are between each known point z_j and the point to be predicted, z_0 .*

Return to Q3 •

Answers

Q4 : *Does the kriging variance depend on the **data values** at the sample points, or on the predicted data value? How can you see this from the equation?* •

A4 : *No and no. There is **no** reference to any **data value** on the right-hand side of the equation.*

*Thus the kriging variance depends **only** on the (1) **sample point configuration** and the (2) **spatial covariance model**.*

Return to Q4 •

Answers

Q5 : *What are the differences between the **predictions** with different block sizes?* •

A5 : *There not much difference, but as the block size increases the extremes are a bit softened; i.e. the hot and cold spots are a bit less pronounced.* *Return to Q5 •*

Answers

Q6 : *What are the differences between the **kriging variances** with different block sizes?* •

A6 : *There is a very large reduction in variance with BK; increasing the block size reduces this still more but not as dramatically.* *Return to Q6* •

Answers

Q7 : *What are the approximate sill and range of the original (blue) and residual (green) variograms?*

The first-order trend surface in this example explained about 40% of the overall variance in the target variable. How is this reflected in the variogram of the residuals from this surface (as shown)? •

A7 : *Sill: 65 CEC-units squared (original) vs. 45 65 CEC-units squared (residual)*

Range: 550 m (original) vs. 350 m (residual).

The trend surface removed much of the total variability, thereby lowering the total sill by about $20/65 = 30\%$. The global trend model acts across the entire area, so that long-range variability due to the trend is removed, hence the shorter range of residuals. *Return to Q7 •*

Answers

Q8 : *What is the relation between the nugget variance of the original (blue) and residual (green) variograms? What should the relation be, in theory?* •

A8 : *The nugget variance is almost the same, approx. 28 CEC-units squared. This agrees with theory: a trend surface can not take out the variability at a single point.*

Technical note: this is because a trend surface is by definition smooth (all partial derivatives are defined).

Return to Q8 •

Answers

Q9 : *What are the major differences in the above figure between OK, global UK, and neighbourhood UK predictions?* •

A9 : *The NW-SE trend somewhat modifies the OK predictions. In particular, notice the lower UK predictions in the area near (180 600, 331 800). This is because the trend surface, which dips towards the SE, is not corrected by nearby observations. Also notice the smaller “halo” of high values in UK near the two high observations on the river bend near (180 200, 330 200). This is because the SE-dipping trend surface predicts low values in this area; the two high points compensate for this in their neighbourhood. In OK they have more influence because there is no contrary trend.*

UK has lower kriging variances than OK, especially at prediction locations with few observations nearby, due to the trend surface.

UK with a local trend is intermediate between global UK and OK in both predictions and variances. In some cases the local trend is even stronger than for UK, for example at the previously-mentioned “cold spot” centred on (180 600, 331 800).

Return to Q9 •

Answers

Q10 : *What are the differences between the **kriging predictions** of OK and CoK? Explain the spatial pattern of these.* •

A10 : *There are fairly large differences between OK and CK in under-sampled areas where the co-variable provides information. For example, near (331200, 179200) (along the river) the under-sampled OK gave under-predictions because of the exclusion of high metal and OM points along the river; when these are added back in the CoK prediction is higher by about 0.3 units $\log(\text{Pb mg kg}^{-1})$. The reverse is true near (332000, 179600).* *Return to Q10 •*

Answers

Q11 : *What are the differences between the **kriging variances** of OK and CoK? Explain the spatial pattern of these.* •

A11 : *Prediction variances are in general larger with CK, with this co-variable; reason: imperfect feature-space correlation as shown also by high nugget in the variogram for OM.*

However, prediction variances are lower in under-sampled areas where a covarible point is available, because the extra information from the co-variable reduces uncertainty. For example, at the extreme SE of the area.

Return to Q11 •