

(1959) had published a paper for the estimation of the normal distribution from censored data. Later work was performed for other distributions such as the 3 parameter lognormal distribution (Cohen, 1976). The first papers concentrated mainly on right censored (survival) data. In Helsel and Cohn (1988) left censored water quality data were analyzed. Despite recent works on the subject such as Shumway et al. (2002) the statistical treatment of censored environmental data is far less applied as it could and should be Helsel (2005). less frequently applied ... should be (Helsel, 2005).

While the treatment of censored environmental data from the classical statistical viewpoint is reasonably well developed this is not the case in spatial statistics. Spatial mapping of variables with censored data is also of great interest and practical importance.

Recently Sedda et al. (2010) presented a methodology to reflect censored data using a simulation approach. In Saito and Goovaerts (2000) the authors addressed the problem of censored and highly skewed variables, and showed that the indicator approach outperforms other geostatistical methods of interpolation.

Variables with non-detects are usually highly skewed which makes their interpolation even more difficult. The high skew of the distributions often leads to problems with the variogram or covariance function estimation. A few large values dominate the experimental curve, and outliers can lead to useless variograms. This problem is partly overcome by the use of indicator variables. However this approach suffers from other deficiencies as demonstrated in this paper.

Purpose of this paper is to develop a methodology to estimate spatial dependence structure from a mixed dataset containing differently censored data. The approach requires as a first step the estimation of the univariate distribution function of the variable under consideration. For this purpose a maximum likelihood method is used. In the next step the spatial dependence is described here with the help of copulas, and the copula parameters are estimated using a maximum likelihood method. After this, the estimated dependence structure is used for the interpolation.

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

◀ ▶

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



The methodology is demonstrated using different water quality parameters obtained from large scale measurement campaigns in South-West Germany. Two highly censored parameters, namely arsenic and deethylatrazin are considered. In order to test the methodology a parameter with no censored data (chloride) is selected and subsequently artificially censored. The methodology is compared to ordinary and indicator kriging using different performance measures.

2 Methodology

2.1 Marginal distribution

Assume that there are n_d measurements with values below the detection limit d_i (note that the detection limits might differ), and for n_z observations a measurement value z_j is given. The empirical distribution function of such observations can only be calculated for values above the largest detection limit. Due to the censoring the mean and the standard deviation cannot be calculated directly, thus the estimation of the parameters θ of a selected parametric distribution via method of moments is not possible. Instead a maximum likelihood method is required. Here one has two choices:

1. To assume a parametric distribution function over the whole domain, and to assess the parameters via maximum likelihood
2. To assume a mixed distribution: for values below a threshold a parametric form is assumed and ~~when above the~~ empirical or a non parametric distribution is considered.

While the first approach is more or less straightforward, it has a few shortcomings. One of them is that outliers might have a very important influence on the parameters of the distribution; the other is that the underlying distribution could be bimodal.

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



In the first case, the

The estimation of the distribution parameters θ can be done using the likelihood function:

$$L_{\text{low}}(\dots, d_i, \dots, z_j \dots | \theta) = \prod_{i=1}^{n_d} F(d_i | \theta) \prod_{j=1}^{n_z} f(z_j | \theta) \quad (1)$$

applied to those values below the local detection limit

where $F(\cdot | \theta)$ is the distribution function and $f(\cdot | \theta)$ the corresponding density with parameter θ .

In the case of the mixed approach we assume that the values below a given threshold z_{lim} follow a parametric distribution, while above the empirical distribution should be considered. Thus the estimation is restricted to those which are below z_{lim} :

$$L_{\text{lim}}(\dots, d_i, \dots, z_j \dots | \theta) = \frac{1}{F(z_{\text{lim}} | \theta)} \prod_{i=1}^{n_d} F(d_i | \theta) \prod_{j=1}^{n_z} f(z_j | \theta) \quad (2)$$

remove second product?

In both cases the logarithm of the likelihood function can be maximized. Above the z_{lim} value a distribution $F_{\text{lim}}(z)$ is assumed.

$$F_{\text{lim}}(z) = \frac{1}{n_{\text{lim}} + 1} \sum_{i=1}^n 1_{z_{\text{lim}} < z < z_i} \quad (3)$$

where n_{lim} is the number of z_i greater than z_{lim} .

The overall distribution function is:

$$G(z) = \begin{cases} F(z | \theta) & \text{if } z \leq z_{\text{lim}} \\ F(z_{\text{lim}} | \theta) + (1 - F(z_{\text{lim}} | \theta)) F_{\text{lim}}(z) & \text{if } z > z_{\text{lim}} \end{cases} \quad (4)$$

Note that the limit z_{lim} is not estimated, but selected as a reasonable limit which is certainly below possible outliers.

outlier observations.

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

This approach allows us to investigate the new variable $U(x) = F(Z(x))$ which has a uniform marginal distribution.

Two copula models, the Gaussian (normal) and the v-transformed normal copula, are considered. The Gaussian copula is described by its correlation matrix Γ .

5 The v-transformed normal copula is parametrized by the transformation parameters m, k and the correlation matrix Γ , which is likely to differ from the Gaussian one.

The v-transformed copula is defined using \mathbf{Y} being an n dimensional normal random variable with $\mathbf{0}^T = (0, \dots, 0)$ mean and Γ correlation matrix $(N(\mathbf{0}, \Gamma))$. All marginals are supposed to have unit variance. Let \mathbf{X} be defined for each coordinate $j = 1, \dots, n$ as:

$$X_j = \begin{cases} k(Y_j - m) & \text{if } Y_j \geq m \\ m - Y_j & \text{if } Y_j < m \end{cases} \quad (7)$$

where k is a positive constants and m is an arbitrary real number. When $k = 1$ this transformation leads to the multivariate non centered χ -square distribution. All one dimensional marginals of \mathbf{X} are identical and have the same distribution function.

10 The parameters of the spatial copula are estimated using the maximum likelihood method.

For the Gaussian copula, as a consequence of the stationarity assumption, the correlations between any two points can be written as a function of the separating vector \mathbf{h} . Then for any set of observations x_1, \dots, x_n the correlation matrix Γ can be written as:

$$\Gamma = \left((\rho_{i,j})_{i,j}^{n,n} \right) \quad (8)$$

where $\rho_{i,j}$ only depends on the vector \mathbf{h} separating the points x_i and x_j :

$$\rho_{i,j} = R(x_i - x_j) = R(\mathbf{h}_{i,j}) \quad (9)$$

For the estimation, the observed values are transformed to the standard normal distribution using:

$$y_k = \Phi_1^{-1}(F(z(x_k))) \quad k = 1, \dots, n_z \quad (10)$$

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



/2 ?? wasn't the median of the interval suggested somewhere in the paper?

$$y_j^d = \Phi_1^{-1}(F(d(x_j))) \quad j = 1, \dots, n_d \quad (11)$$

, rather like a variogram

Here $\Phi_1(\cdot)$ is the distribution function of the standard normal distribution $N(0,1)$.

The variable y is now normal with data below the detection limit denoted by y_j^d . The correlation function $R(\cdot, \beta)$ is assumed to have a parametric form with the parameter vector β . The likelihood function in this case can be written as:

$$L(\beta) = \prod_{(j,k) \in I_1} \phi_2(y_j, y_k, R(\mathbf{h}_{j,k}, \beta)) \prod_{(j,k) \in I_2} \Phi_1\left(\frac{y_j^d - y_k R(\mathbf{h}_{j,k}, \beta)}{\sqrt{1 - R(\mathbf{h}_{j,k}, \beta)^2}}\right) \prod_{(j,k) \in I_3} \Phi_2(y_j^d, y_k^d, R(\mathbf{h}_{j,k}, \beta)) \quad (12)$$

Here $\Phi_2(x, y, r)$ is ~~the distribution function of~~ the 2 dimensional normal distribution with correlation r and standard normal marginal distributions $N(0,1)$ and $\phi_2(x, y, r)$ is its density function. The set I_1 contains pairs of locations with both variables being measured exactly. In I_2 pairs are listed which consist of an exact observation and a **below detection limit value**. Finally, I_3 contains pairs with values below the detection limit. The logarithm of the likelihood function is maximized numerically.

The above procedure might require a lot of computation effort if the number of observations is large. Instead one can reduce the number of pairs considered in Eq. (12) by selecting different distance classes and taking each observation exactly M times as a member of a pair. This way one can avoid clustering effects. **That's a smart idea.**

A similar but slightly more complicated procedure has to be used for the estimation of the parameters of the v-copula. In this case the variable Z is first transformed to:

$$y_k = H_1^{-1}(F(z(x_k))) \quad k = 1, \dots, n_z \quad (13)$$

value below d_j .

$$y_j^d = H_1^{-1}(F(d(x_j))) \quad j = 1, \dots, n_d \quad (14)$$

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



Here $H_1(\cdot)$ is the univariate distribution function of the v-transformed normal distribution. This can be written as:

$$H_1(y) = \Phi_1\left(\left(\frac{y}{k}\right) + m\right) - \Phi_1(m - y) \quad (15)$$

the copula density function is illustrated in the figure on page 5283, which might help the reader ...

The likelihood function in this case is:

$$L(\beta) = \prod_{(j,k) \in I_1} h_2(y_j, y_k, \beta) \prod_{(j,k) \in I_2} H_c(y_j^d, y_k, \beta) \prod_{(j,k) \in I_3} H_2(y_j^d, y_k^d, \beta) \quad (16)$$

The sets I_1, I_2 and I_3 are defined as for the Gaussian case. $H_2(\dots)$ is the distribution function of the bivariate v-transformed distribution:

$$H_2(y_1, y_2, \beta) = \Phi_2\left(\left(\frac{y_1}{k}\right) + m, \left(\frac{y_2}{k}\right) + m, R(\mathbf{h}_{j,k}, \beta)\right) + \Phi_2(m - y_1, m - y_2, R(\mathbf{h}_{j,k}, \beta)) - \Phi_2\left(\left(\frac{y_1}{k}\right) + m, m - y_2, R(\mathbf{h}_{j,k}, \beta)\right) - \Phi_2\left(m - y_1, \left(\frac{y_2}{k}\right) + m, R(\mathbf{h}_{j,k}, \beta)\right) \quad (17)$$

change j,k to 1,2 for consistency?

5 Here $R(\mathbf{h}_{j,k}, \beta)$ is the correlation function of the Gaussian variable Y and $h_2(\dots)$ is the density function corresponding to H_2 . The bivariate function $H_c(\dots)$ is obtained via integration of the density:

$$H_c(y_1, y_2, \beta) = \int_{-\infty}^{y_1} h(y, y_2) dy \quad (18)$$

mixed

i don't understand this phrase - please expand

As the density h is a weighted sum of normal densities, the corresponding integral can be calculated for each term separately, which is similar to the normal case.

10 Due to the complex form of the overall likelihood function, a numerical optimization of the log-likelihood function is done.

complicated

Different forms of the correlation function can be considered – such as the exponential:

$$R(\mathbf{h}, A, B) = \begin{cases} 0 & \text{if } \|\mathbf{h}\| = 0 \\ B \exp\left(-\frac{\|\mathbf{h}\|}{A}\right) & \text{if } \|\mathbf{h}\| > 0 \end{cases} \quad (19)$$

single bar '|' for modulus?

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



where $0 \leq B \leq 1$ and $A > 0$.

3 Interpolation

Once the parameters of the correlation function (A, B) and for the v-transformed copula the parameters of the v-transformation (m, k) are estimated the interpolation can be carried out. In order to reduce the complexity of the problem, interpolation will be done using a limited number of neighboring observations. Due to an umbrella effect similar to as for ordinary kriging, observations which are behind other observations have a minor influence on the conditional distribution. Further, this restriction to local neighborhoods relaxes the assumption of stationarity to a kind of local stationarity. An example in Bárdossy and Li (2008) demonstrates that this assumption does not significantly alter the results of interpolation.

The goal of interpolation is to find the density of the random variable $Z(x)$ conditioned on the available censored and uncensored observations. The conditional density $f_x(z)$ for location x can be written as:

$$f_x(z) = P(Z(x) = z | Z(x_i) < d_i, i = 1, \dots, n_d; Z(x_j) = z_j, j = 1, \dots, n_z) = \frac{P(Z(x) = z, Z(x_i) < d_i, i = 1, \dots, n_d | Z(x_j) = z_j, j = 1, \dots, n_z)}{P(Z(x_i) < d_i, i = 1, \dots, n_d | Z(x_j) = z_j, j = 1, \dots, n_z)} = \frac{P(Z(x_i) < d_i, i = 1, \dots, n_d | Z(x) = z, Z(x_j) = z_j, j = 1, \dots, n_z) P(Z(x) = z)}{P(Z(x_i) < d_i, i = 1, \dots, n_d | Z(x_j) = z_j, j = 1, \dots, n_z)} \quad (20)$$

Both the **nominator** and the denominator of the last expression are conditional multivariate distribution function values.

numerator (there are other instances)

it would make it easier to read if commas were introduced here

values which require integration of the multinormal in n_d dimensions.

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



For the normal copula case, Eq. (20) can be written with the help of the transformed variable Y

it would make it easier to read if commas (or semicolons as in equation (24) below) were introduced here

$$f_x(z) = \frac{P\left(Y(x_i) < y_i^d; i = 1, \dots, n_d | Y(x) = y; Y(x_j) = y_j; j = 1, \dots, n_z\right) P(Y(x) = y)}{P\left(Y(x_i) < y_i^d; i = 1, \dots, n_d | Y(x_j) = y_j; j = 1, \dots, n_z\right)} \quad (21)$$

The conditional distribution of a multivariate normal distribution is itself multivariate normal with expectation μ_c^0 and covariance matrix Γ_c^0 with:

$$\Gamma_c^0 = \Gamma_{00} - \Gamma_{01} \Gamma_{11}^{-1} \Gamma_{01}^T \quad (22)$$

The expected value of the conditional is:

$$\mu_c^0 = \Gamma_{01} \Gamma_{11}^{-1} \mathbf{y} \quad (23)$$

$\mathbf{y}^T = (y, y_1, \dots, y_{n_z})$. The matrices Γ_{00} , Γ_{01} and Γ_{11} are the correlation matrices corresponding to the pairs of observations with censored and uncensored data, calculated with the correlation function $R(h)$.

Thus for the conditional probability in the nominator in (20) can be calculated as:

$$P\left(Z(x_j^0) < d_j; j = 1, \dots, n_d | Z(x) = z; Z(x_j^1) = z_j; j = 1, \dots, n_z\right) = \Phi_{\mu_c^0, \Gamma_c^0}(y, y_1, \dots, y_{n_z}) \quad (24)$$

where $\Phi_{\mu_c^0, \Gamma_c^0}$ is the distribution function of $N(\mu_c^0, \Gamma_c^0)$. Values of the multivariate normal distribution function can be calculated by numerical integration, for example using Genz and Bretz (2002). The denominator in (20) requires the same type of calculations.

The denominator is independent of the value z and can be calculated exactly as the nominator. Note that the point for which the interpolation has to be carried out is considered as a pseudo observation with the observed value z . Thus the nominator has to be evaluated for a number of possible z values to estimate the conditional density.

For the v -transformed copula the interpolation procedure is slightly more difficult, but as the n -dimensional density of the v -transformed variable is a weighted sum of 2^n normal densities the calculation procedure is similar. However, we will not go into further details here.

4 Application and results

The above described methodology was applied to a regional groundwater pollution investigation. Two censored variables and an artificially censored variables were used to demonstrate the methods, and to compare them to traditional interpolations.

4.1 Investigation area

An extensive dataset consisting of more than 2500 measurements of groundwater quality parameters of the near surface groundwater layer in Baden-Württemberg were used to illustrate the methodology. Three quality parameters namely deethylatrazine – degradation product of atrazine – arsenic and chloride were selected for this study.

While the first two parameters are heavily censored the chloride exceed the detection limit in 99.9% of the cases. This variable is artificially censored using different thresholds in order to show the effectiveness of the method.

Table 1 shows the basic statistics for the selected data. Note the high positive skewness for all variables. This alone would leave to substantial difficulties in estimating spatial correlation functions, even in the case if most values had been above the detection limit.

4.2 Parameter estimation

As a first step the marginal distributions were estimated using the approach described in Sect. 2.1. Figure 1 shows the distribution functions for arsenic. The estimation method was compared to the full maximum likelihood (which would correspond to

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



$z_{lim} > \max(z_j, j = 1, \dots, I)$). One can see that the traditional maximum likelihood estimation is strongly influenced by outliers, leading to unrealistic, and unacceptable results. In contrast, setting z_{lim} such that $z_{lim} > \max(d_j, j = 1, \dots, J)$ and bearing in mind that there are at least a few (30 or more) z_j values below z_{lim} , leads to a good fit of the observed values. As a rule of thumb the value of z_{lim} was artificially chosen 50 % above the largest detection limit.

In order to investigate the quality of the extension of the distribution to low values the observed chloride concentration values were artificially censored. Detection limits were set to the 15, 25, 35, 45, 55, 65, 75 and 85 % value of the distribution. Figure 2 shows distribution functions corresponding to different detection limits for chloride. Note that in order to see any differences the x-axis is shown on a logarithmic scale. All distribution functions are very similar, showing that the upper middle part of the distribution can be well used to extend it to low values.

The parameters of the spatial structure were estimated both for a normal and a v-transformed normal copula. An exponential spatial correlation function was assumed. Table 2 shows the parameters of the spatial copulas for the selected variables. The copula fits are very different. While for arsenic the correlation function of the normal copula has a high B value indicating a strong spatial structure, for the v-transformed copula the B is much lower. For deethylatrazine the situation is inverted: the v-transformed copula shows a strong spatial link and the normal nearly no spatial correlations.

4.3 Interpolation

In order to illustrate the properties of the interpolation method illustrative examples are first considered. Assume that the value at the center of a square is to be estimated, with observations at the four corners:

1. Assume all four corners have exact values.

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



4.4 Comparison with other interpolation methods

As an alternative ordinary kriging (OK) was used for interpolation. Three different treatments of the values below the detection limit were considered:

1. All values below the detection limit were set to zero
- 5 2. All values below the detection limit were set to the half of the corresponding detection limit
3. All values below the detection limit were set to the corresponding detection limit

Empirical variograms were calculated for each case. Additionally the empirical variogram was calculated from the exact values only. Figure 6 shows the graph of these variograms for deethylatrazine. The exact values lead to a variogram without any structure and with the highest variance. The datasets with replaced values show a much lower variability and the replacement with zeros increases the variability only very slightly. These variograms do not show a spatial structure. Only after the removal of a few extremes, which were considered as outliers one could obtain a reasonable variogram. This example gives a good idea about the difficulties involved in the assessment of a reasonable variogram. The same procedure was carried out for arsenic and chloride. In the later case the variograms were calculated for different levels of censoring. A cross validation using OK was performed for each parameter and each censoring.

Another popular method to treat highly skewed variables is indicator kriging (IK). The indicator corresponding to a cutoff value α is defined as:

$$I_{\alpha}(Z(x)) = \begin{cases} 0 & \text{if } Z(x) > \alpha \\ 1 & \text{if } Z(x) \leq \alpha \end{cases} \quad (25)$$

20 Indicator variograms are calculated for a set of α values. These do not suffer from the problem of outliers. A subsequent IK leads for each x and α to an estimated

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



spatial mean is below the 55% value of the distribution. Thus for the high levels of censoring the interpolated mean is below the lowest measured value.

As a next step for all three variables and all interpolation methods a cross validation was carried out. The evaluation of the cross validation results is not straightforward due to the censoring. The usual squared error is, even for the exact values, not appropriate as the distributions are highly skewed, and some extreme outliers would dominate this measure. Instead this measure was calculated by leaving out the upper 1% of the measured values, ensuring that outliers were not considered for the calculation. Further the rank correlation for the exact values was calculated. Additionally the LEPS score Ward and Folland (1991) was calculated to evaluate the fit in the probability space.

$$\text{LEPS} = \frac{1}{n} \sum_{i=1}^n |G_z(z(\mathbf{x}_i)) - G_z(z^*(\mathbf{x}_i))| \quad (26)$$

For the measurements below the detection limit the average of the probabilities to be below the detection limit was calculated.

5 Results for the two censored variables and for an artificially censored case chloride are displayed in Tables 3 and 4. As one can see the copula based approaches outperform the ordinary and the indicator kriging. Note that the mean squared error, the rank correlation and the LEPS score were all calculated for the exact measurements only. From the two copula models the v-copula allowing a non-symmetrical dependence is slightly better than the Gaussian.

10 For the artificially censored mean squared error, rank correlation and LEPS score were calculated using all data without considering the artificial censoring. Thus these measures represent a realistic measure of interpolation quality. The results are shown in Table 5. Note that ordinary kriging has a very high mean squared error. This is caused by the high skewness of the marginal distribution which had much less influence on the indicator and copula approaches.

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



For interpolation and for possible random simulation of the fields a good measure of uncertainty is of great importance. As the kriging variance is only data configuration but not data value dependent (especially for skewed distributions c.f. Journel, 1988) it is not a good measure of uncertainty. The indicator approach provides estimates of the local conditional distribution functions. As it is not directly considering the estimation uncertainty (all indicator values are interpolated values with no uncertainty associated) it does not provide a good uncertainty measure. The copula approach yields full probability distributions for each location, thus arbitrary confidence intervals can be derived. Figure 11 shows the width of the 80 % confidence interval obtained using v-copula based interpolation and the kriging standard deviations from OK for deethylatrazin. One can see that the estimation quality of the copula based interpolation is very heterogeneous over the whole domain. Regions with high observed values the confidence intervals are wide, in low areas narrow. For ordinary kriging the estimation error (kriging standard deviation) is small close to points with measured values, irrespective of the observed values.

In order to validate the confidence intervals the frequency of observations within the 80 % confidence interval (obtained from cross validation) was calculated. Figure 9 shows the percentage of chloride values falling into the 80 % confidence interval for different censoring levels obtained using the v-copula and the Gauss copula. As one can see for the v-copula the frequency is close to the target 80 % for all censoring levels while for the Gauss copula the confidence intervals become meaningless above 35 % censoring.

5 Conclusions

In this paper a methodology for the interpolation of variables with data below a detection limit was developed. As a first step the marginal distributions were estimated using a mixed approach which entailed a maximum likelihood method for the lower values and the empirical distribution for the high values. This procedure provides a robust

why not put copula interpolation CL to 56% to match the 1 stdev CL for OK?

estimator for the low concentrations without the negative influence of possible outliers. Using the fitted distributions the variables were transformed to the unit interval and their spatial copula was assessed, assuming spatial stationarity. Values below the detection limit are considered in a maximum likelihood estimation of the spatial copula parameters. Interpolation was done by calculating the conditional distributions for each location. The conditions include both the measurements as exact values and the below detection limit observations as inequality constraints.

The copula based interpolation is exact at the observation locations; the interpolated value equals the observed value. For locations with censored observations the method provides an updated distribution function which differs from the constrained marginal. Other procedures such as indicator kriging with inequality constraints do not update distributions at observation locations.

Investigations based on the artificially censored dataset show that the copula-based approaches remain unbiased even for large degrees of censoring. Among the kriging approaches only ordinary kriging with setting the censored values equal to the half of the corresponding detection limit did not show a systematic error for higher detection limits. This choice is clearly better than setting the values below the detection limit equal to the detection limit, or setting them all equal to zero, which both lead to systematic errors. Indicator kriging also shows a systematic bias increasing with the detection limit.

The copula-based approaches outperform ordinary and indicator kriging in their interpolation accuracy. Indicator kriging is only slightly worse than the copula based interpolation, while ordinary kriging with all different considerations of the values below detection limit are the poorest estimators.

The main advantage of the copula based approaches is in the estimation of the interpolation uncertainty. While ordinary kriging yields unrealistic estimation variances depending only on the configuration of the measurement locations, the copula-based interpolation yields reasonable confidence intervals. The v-copula based approach yields more realistic confidence intervals than the Gaussian alternative.

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



The suggested approach can be extended to handle any kind of inequality constraints both for spatial structure assessment and for interpolation.

The model can serve as a basis for conditional spatial simulation. It is imaginable to extend the model to a Bayesian approach where prior distributions are assigned to individual locations.

Acknowledgements. Research leading to this paper was supported by the German Science Foundation (DFG), project number Ba-1150/12-2.

References

- Bárdossy, A.: Copula-based geostatistical models for groundwater quality parameters, *Water Resour. Res.*, 42, W11416, doi:10.1029/2005WR004754, 2006. 5268
- Bárdossy, A. and Li, J.: Geostatistical interpolation using copulas, *Water Resour. Res.*, 44, W07412, doi:10.1029/2007WR006115, 2008. 5268, 5272
- Cohen, C.: Simplified Estimators for the Normal Distribution When Samples Are Singly Censored or Truncated, *Technometrics*, 1, 217–237, 1959. 5264
- Cohen, C.: Progressively Censored Sampling in the Three Parameter Log-Normal Distribution, *Technometrics*, 18, 99–103, 1976. 5265
- Genz, A. and Bretz, F.: Comparison of Methods for the Computation of Multivariate t-Probabilities, *J. Comp. Graph. Stat.*, 11, 950–971, 2002. 5273
- Helsel, D. R.: More than obvious: Better methods for interpreting nondetect data., *Environ. Sci. Technol.*, 39, 419A–423A, 2005. 5265
- Helsel, D. R. and Cohn, T. A.: Estimation of descriptive statistics for multiply censored water quality data, *Water Resour. Res.*, 24, 1997–2004, 1988. 5265
- Journal, A. G.: New Distance Measures: The Route Toward Truly Non-Gaussian Geostatistics, *Math. Geol.*, 20, 459–475, 1988. 5280
- Roth, C.: Is lognormal kriging suitable for local estimation?, *Math. Geol.*, 30, 999-1009, 1998. 5278
- Saito, H. and Goovaerts, P.: Geostatistical interpolation of positively skewed and censored data in a dioxin-contaminated site, *Environ. Sci. Technol.*, 44, 4228–4235, 2000. 5265

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

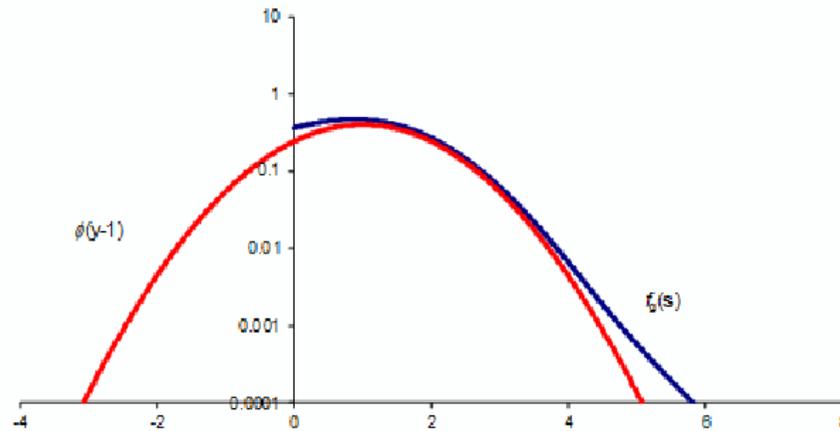
Full Screen / Esc

Printer-friendly Version

Interactive Discussion



- Sedda, L., Atkinson, P. M., Barca, E., and Passarella, G.: Imputing censored data with desirable spatial covariance function properties using simulated annealing, *J. Geogr. Syst.*, 36, 3345–3353, 2010. 5265
- Shumway, R., Azari, R., and Kayhanian, M.: Statistical Approaches to Estimating Mean Water Quality Concentrations with Detection Limits, *Environ. Sci. Technol.*, 36, 3345–3353, 2002. 5265
- Ward, M. and Folland, C.: Prediction of seasonal rainfall in the Nordeste of Brazil using eigenvectors of sea-surface temperature, *International Journal Climatology*, 11, 711–743, 1991. 5279



The pdf of the $N(1,1)$ distribution compared with that of the v -copula with $m = 1$. The log-scale is used to emphasise the tail behaviour.

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



Table 1. Basic statistics of the investigated variables mean, standard deviation and skewness are calculated from values above the detection limit.

	Number of observations	Number of above DL	Statistics of values > Detection limit			
			Mean	Standard deviation	Skewness	Maximum
Arsenic	2234	979	0.002733	0.007392	13.4	0.1618
deethylatrazine	2848	403	0.064243	0.068316	4.5	0.68
Chloride	2805	2801	39.9	165.8	30.3	6940.0

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

|◀ ▶|

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



Table 2. Parameters of the fitted copulas.

	Gauss copula		V-transformed copula			
	<i>B</i>	<i>A</i>	<i>B</i>	<i>A</i>	<i>m</i>	<i>k</i>
Arsenic	0.750	1325	0.810	49000	1.78	0.376
deethylatrazine	0.030	669	0.579	35000	0.29	2.469
Chloride	0.620	11539	0.449	27500	1.98	0.147

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD
8, 5263–5299, 2011

Interpolation of groundwater quality parameters
A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

|◀ ▶|

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



Table 3. Cross validation results for Arsenic.

Measure	V-copula	Gauss-copula	Indicator Kriging	Ordinary Kriging 50 % of Detection limit
MSQE	3.7×10^{-6}	1.0×10^{-5}	5.3×10^{-5}	1.0×10^{-5}
Rank correlation	0.32	0.32	0.33	0.33
LEPS Score	0.142	0.154	0.142	0.159
Mean probability for < DTL	0.610	0.559	0.042	0.437

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

|◀ ▶|

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



Table 4. Cross validation results for deethylatrazin.

Measure	V-copula	Gauss-copula	Indicator Kriging	Ordinary Kriging 50 % of Detection limit
MSQE	5.1×10^{-4}	3.0×10^{-3}	5.0×10^{-3}	1.7×10^{-3}
Rank correlation	0.44	0.31	0.40	0.48
LEPS Score	0.168	0.311	0.100	0.110
Mean probability for < DTL	0.869	0.888	0.560	0.650

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

|◀ ▶|

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



Table 5. Cross validation results for Chloride with 45 % artificial censoring.

Measure	V-copula	Gauss-copula	Indicator Kriging	Ordinary Kriging 50 % of Detection limit
MSQE	273.1	251.8	298.6	2922.5
Rank correlation	0.61	0.61	0.58	0.45
LEPS Score	0.186	0.174	0.191	0.150
Mean probability for < DTL	0.593	0.555	0.000	0.390

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD
8, 5263–5299, 2011

**Interpolation of
groundwater quality
parameters**

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

|◀ ▶|

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



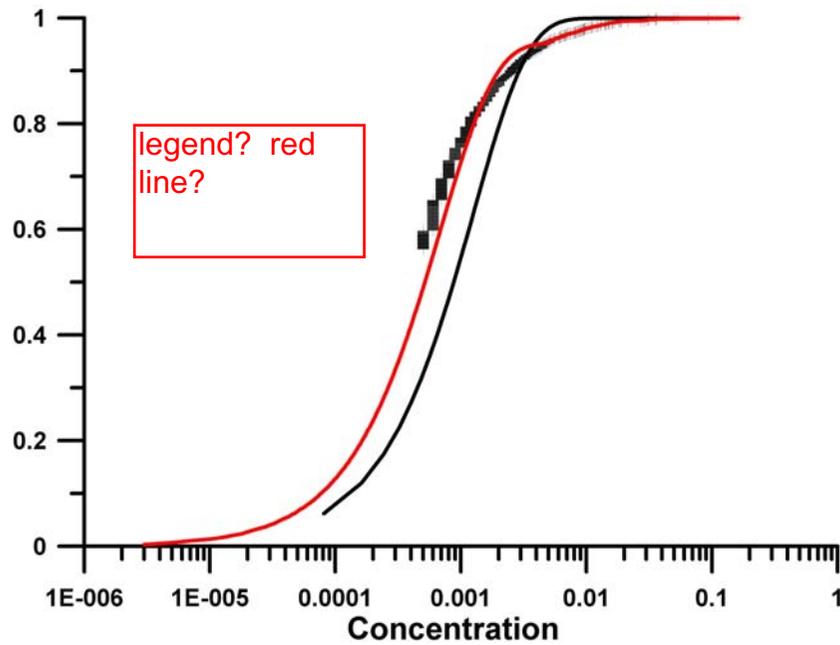


Fig. 1. The distribution of the observed arsenic concentrations and the distributions obtained via maximum likelihood for the whole dataset (black line) and with setting z_{lim} to 1.5 times the highest detection limit.

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

|◀ ▶|

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

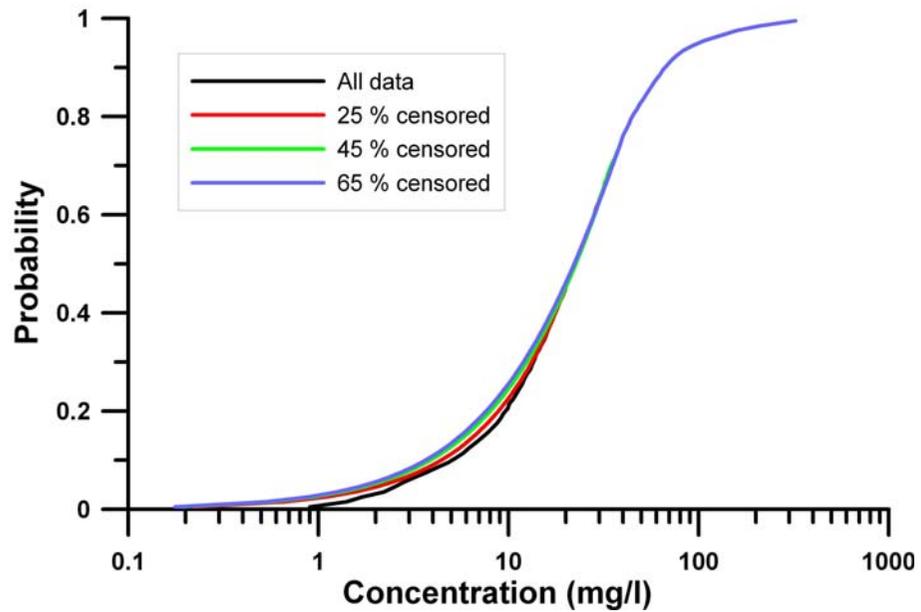


Fig. 2. The distribution of chloride concentrations and the estimated distributions corresponding to different degrees of censoring.

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

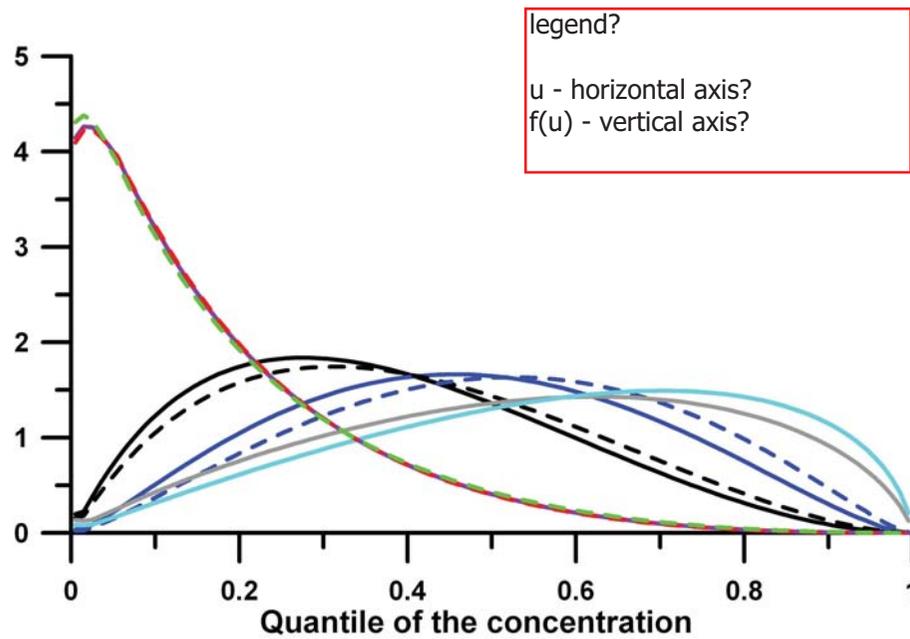
Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion





Conditional copula densities ...

Fig. 3. Conditional densities obtained for the center of a square using different data at the corners. The light blue and the gray lines correspond to exact observed values with the 75 % quantile at each corner and with the 75 % quantile at three corners and the 95 % quantile at the forth. The dark blue lines correspond to one exact value at the 95 % value and three corners below the detection limit which is at the 75 % value (solid) and at the 95 % value (dashed line). The black lines correspond to one exact value at the 75 % value and three corners below the detection limit which is at the 75 % value (solid) and at the 95 % value (dashed line). The red and green lines show results for the case if all values are below the detection limit.

limits of 75% and 95% respectively

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



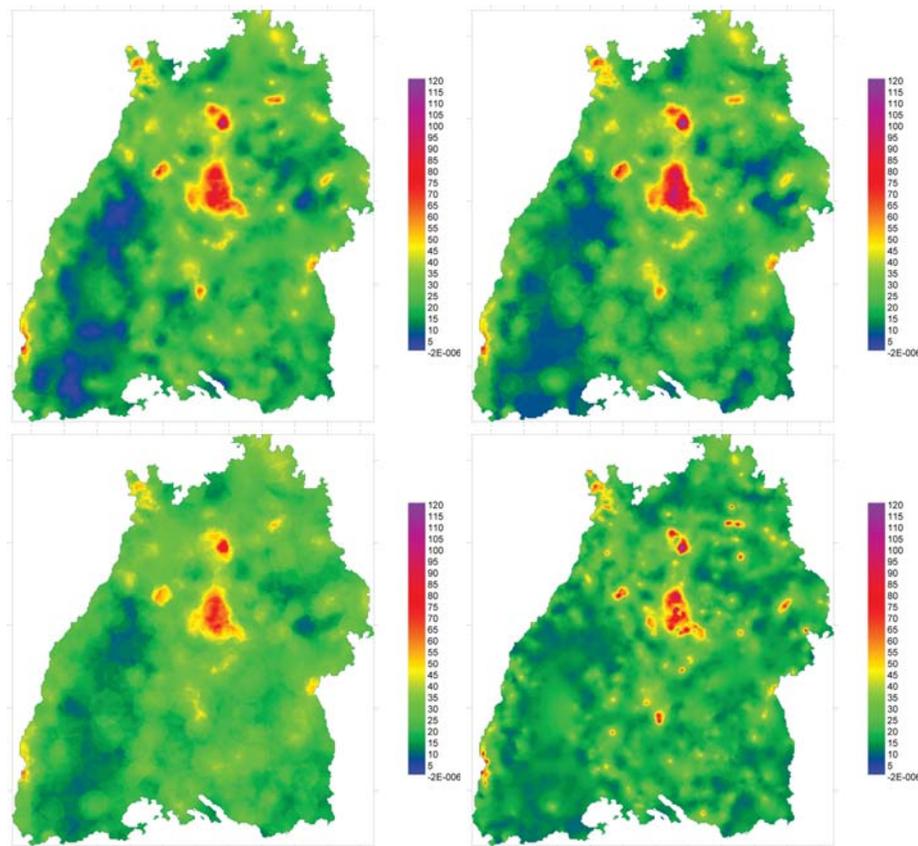


Fig. 4. Interpolated chloride concentrations for different grades of censoring.

legend?

from text: interpolated maps for chloride using all observations and three different maps using 25 %, 45% and 65% censoring.

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Introduction

Conclusions

References

Tables

Figures

◀

▶

◀

▶

Back

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



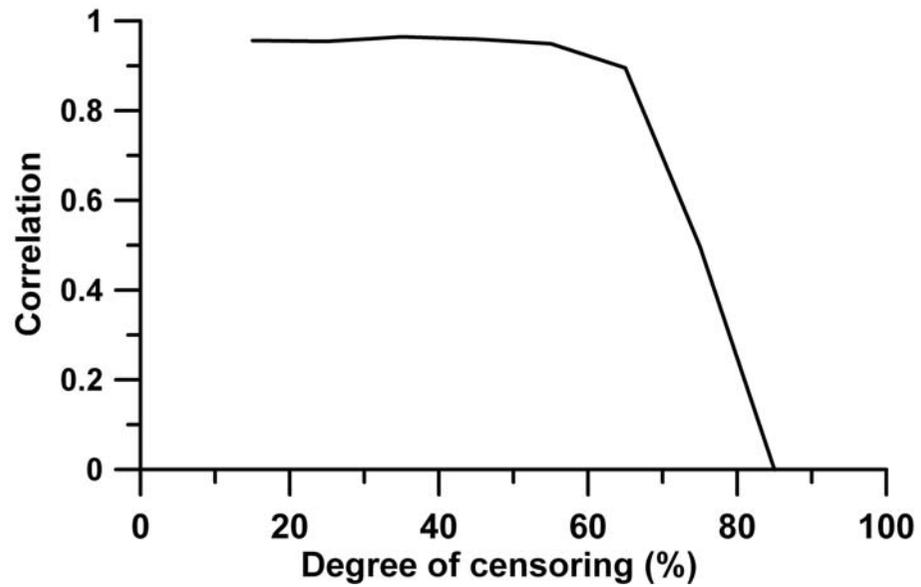


Fig. 5. Correlation between the interpolated map of Chloride and **the** maps interpolated from censored data.

this graph has been calculated from more than three (eight? from the kinks in the line) maps

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

|◀ ▶|

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

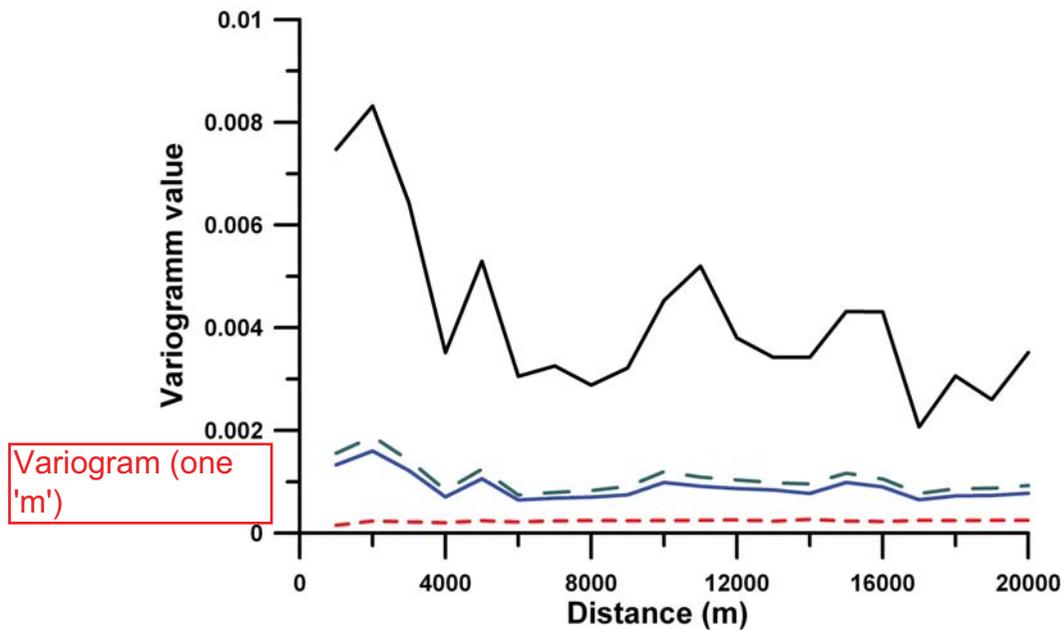


Fig. 6. Empirical variograms calculated for deethylatrasine, using exact data only (black solid), using nondetects replaced by zero (blue dashed) or by the detection limit (blue solid) and using nondetects replaced by zero and removal of outliers.

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract

Conclusions

Tables

◀

▶

Back

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

Introduction

References

Figures

▶

◀

Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

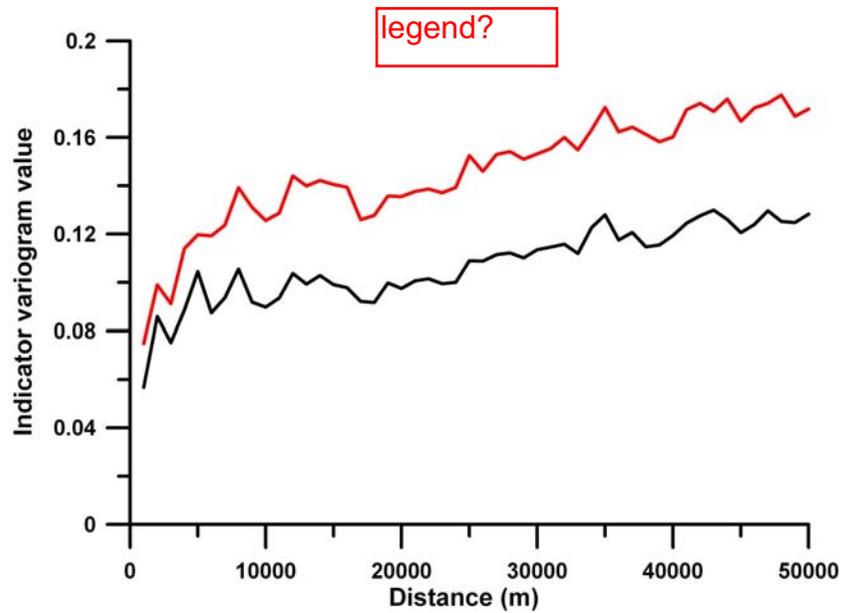


Fig. 7. Empirical indicator variograms calculated for deethylatrasine for the 85% and 90% values of the distribution.

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

|◀ ▶|

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

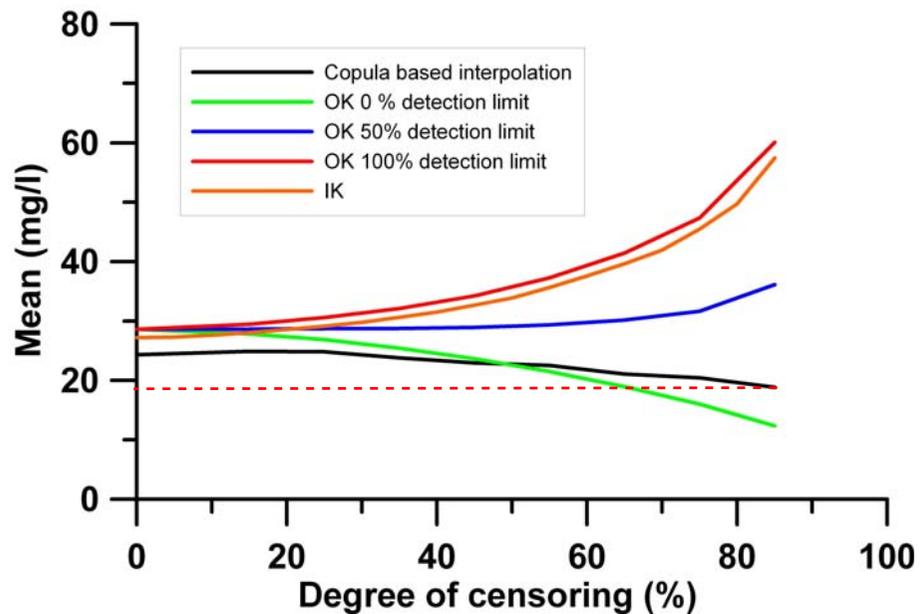


Fig. 8. Mean of the interpolated maps of Chloride for different degrees of censoring and different interpolations.

it is worth noting that the copula-based interpolation starts from a lower mean (25?) than the others and drops only to 19? with 85% decimation

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

◀ ▶

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion

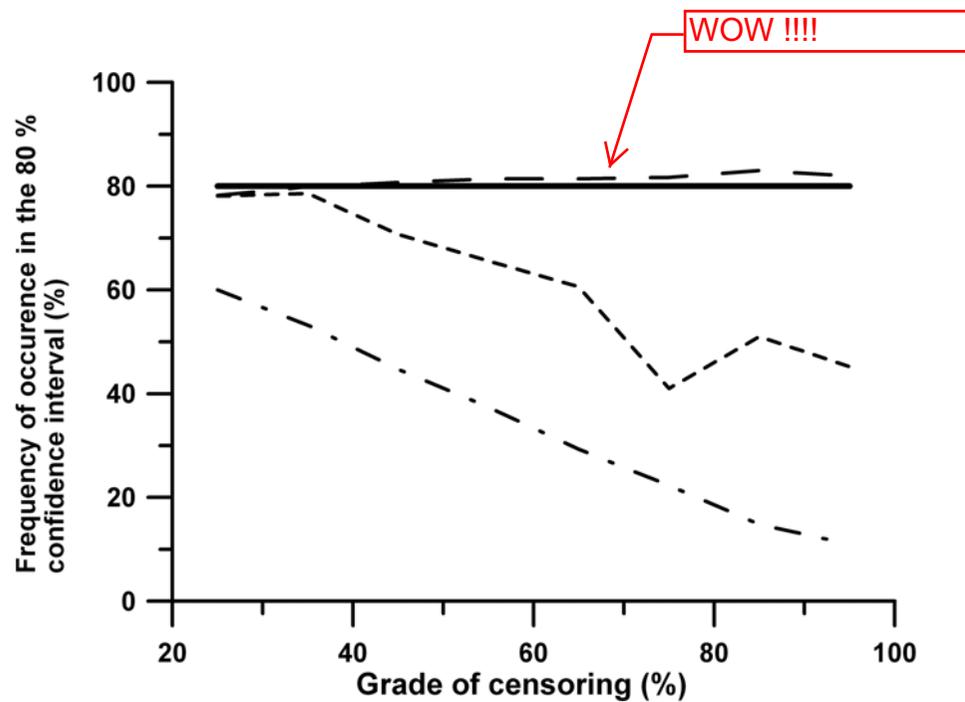


Fig. 9. Frequency of observations in the 80 % confidence interval for V-copula based interpolation (long dashes) and Gauss-copula based interpolation (short dashes) and indicator kriging (dashed dotted line) for different grades of censoring of Chloride.

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

|◀ ▶|

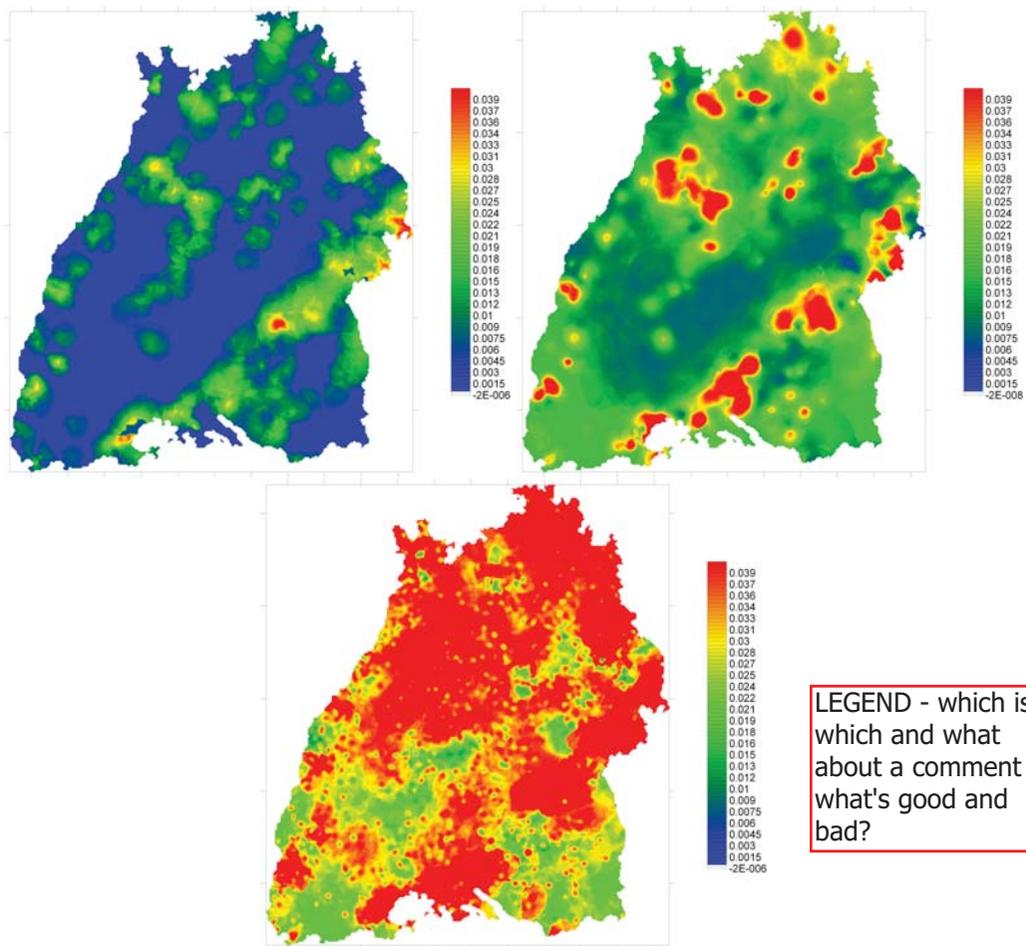
◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



LEGEND - which is which and what about a comment on what's good and bad?

Fig. 10. Interpolated deethylatrazine concentrations using different interpolation methods. For OK the values were set to the detection limit.

not clear what this means

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

◀ ▶

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



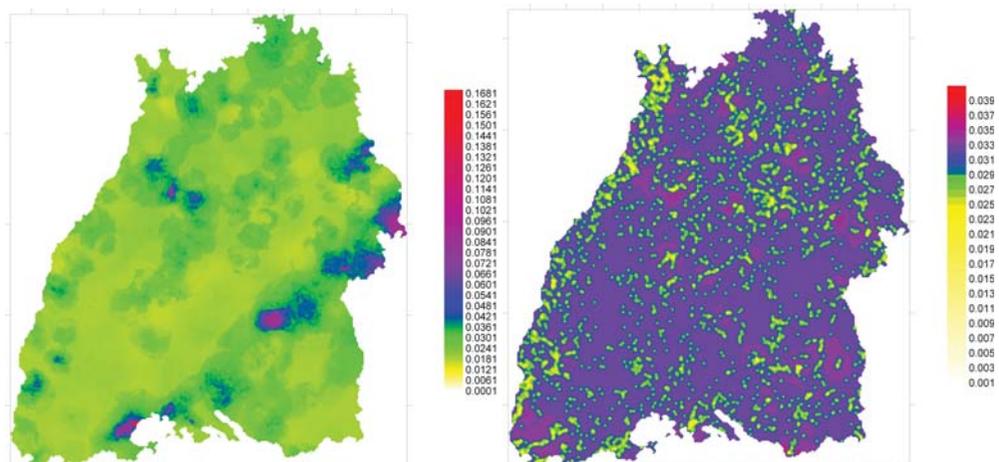


Fig. 11. Uncertainty maps for deethylatrazin: left the length of the 80% confidence interval obtained via v-copula based interpolation, right the kriging standard deviation obtained by OK.

note differences in scale and almost constant size of OK interval.

why not put copula interpolation CL to 56% to match the 1 stdev CL for OK? is it difficult to give them the same colour range on the legends for ease of comparison?

Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper | Discussion Paper

HESSD

8, 5263–5299, 2011

Interpolation of groundwater quality parameters

A. Bárdossy

Title Page

Abstract Introduction

Conclusions References

Tables Figures

◀ ▶

◀ ▶

Back Close

Full Screen / Esc

Printer-friendly Version

Interactive Discussion