

Spatial Error and Interpolation Uncertainty Appraisal Within Geographic Information Systems

¹J. Negreiros, ¹M. Painho, ²M.A. Aguilar and ²F.J. Aguilar

¹Instituto Superior de Estatística e Gestão de Informação-Universidade Nova de Lisboa, Campus de Campolide, 1070-312 Lisboa, Portugal

²Departamento de Ingeniería Rural-Escuela Politécnica Superior, Universidad de Almería, La Cañada de San Urbano S/N 04120 Almería, Spain

Abstract: From a GIS software point of view, spatial error and quality control of spatial data has been almost forgotten with the exception of geostatistics. This review study pretends to re-evaluate major sources of errors in spatial data. This includes qualitative and quantitative accuracy, data quality principles and interpolation uncertainty perceptions. As an illustration purpose, an empirical analysis of 98 samples of Pb lead contaminant is also presented along this study.

Key words: Geographical Information Systems (GIS), data quality, spatial error, uncertainty assessment

INTRODUCTION

Data quality might be defined as the general intent of fitness of a particular dataset for a particular use that one may have in mind for the data (Chrisman, 1992). Surveying skills, support of samples, measurement device accuracy, the choice of map projections and spheroids, rounding errors and visual presentation are some sources of uncertainty. In terms of time process, the loss of spatial data quality may occur in different phases, that is, during digitalization, storage, documentation, analysis, presentation and through the use to which they are put (Fig. 1). Hence, error must not be treated as a potentially embarrassing inconvenience, because error provides a critical component in judging fitness (a measure of how well the geometric representation matches the original spatial representation) for use (Chrisman, 1992).

Still, the major flaw is the absence of GIS tools for error evaluation. GIS users always have in mind that measurement errors and short-range variation contribute to local uncertainty, which can be very large indeed (Burrough, 1993). In particular, GIS error propagation is quite important during an overlay operation of 2 different thematic layers or 2 layers of the same theme but originally at different scales. Certainly, the Monte Carlo simulation can play an important role to model these effects. Tomlinson (1992), substantiate a distinctive belief stating if we are to convert this huge amount of human experience and effort to digital form, we need better

digitizing methods, that is, data selection, error correction, coordinates conversion, editing and reformatting, in some order sequence. Independent from the total volume of spatial data that is considered, it is crucial to know the ratio between images and attribute data. As the image versus attribute ratio increases, the spatial data handling problems increase as well in a hasty way (Calkins, 1992).

Error can be questionable in an estimation process, too. According to Isaaks and Srivastava (1989), it is always be modest and confess ignorance in large unsampled area instead involving extrapolating data values over large distances that, in general, can lead to quite misleading results. Further, the presence of a large quantity of data samples for huge areas might create a problem for the geostatistician since, he/she must be aware, which sampling distribution should be choose for statistical processing in order to represent the whole exhaustive data sample distribution.

Quite often, the quality of GIS software is judged by the visual appearance to the human eye. This is surprising given the costs of data acquisition and the investments that are linked to the use of GIS (Burrough and MacDonnell, 1998). In general, these cost round 60% of the total project budget. Only in the field of geostatistics, uncertainty and error measurement has been developed to deal with both issues, particular with Indicator Kriging (IK), Probability Kriging (PK), indicator simulation, cross-validation and Kriging variance.

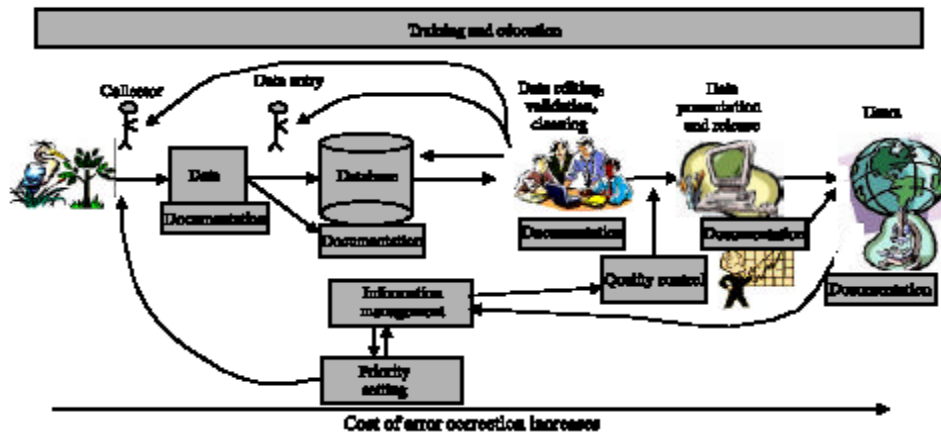


Fig. 1: Information management Chain showing that the cost of error correction increases as one progresses along the chain. Good documentation, education and training are integral to all steps (Chapman, 2008)

To review errors, quality control and uncertainty assessment in a GIS framework becomes the core of this state-of-the-art study.

SOURCES OF GIS ERRORS

The capability of an institution to produce high quality data is influenced by the level of staff experience and laboratory facilities. Yet, to determine spatial data quality is not an easy task. A major suggestion by Chapman (2008) regarding principles of data quality is shown in Table 1 where all data must be attached by a documentation quality field. This will provide an indication of the certainty of the identification.

Chrisman (1992) argued that error is an integral part of spatial information processing and it should be recognized as a fundamental dimensional issue of spatial data. Hence, GIS should put a significant effort into the development of methods to report and visualize databases error like outliers, logical consistency values, missing fields, typographic errors or non-atomic data values (more than one fact entered into a single field). Isaaks and Srivastava (1989), suggest to sort the data and to examine the extreme values, to locate them on a map, especially if they are isolated, or to check coordinates errors in order to produce clean data. According to both authors, extreme values should be deleted from the data set on estimated studies.

Painho (1992) and confirms Chrisman (1992) view stating that error should be added to position and attribute GIS data. By directly recognizing error, it may be possible to confine, it to acceptable limits. It seems, it is more efficient to achieve a level of error compatible with the purpose of the analyses that are to be performed than swept away under the carpet. Still, error cannot always be avoided cheaply or easily.

Table 1: Possible classification of documentation quality

1.	Identified by world expert in the taxa with high certainty
2.	Identified by world expert in the taxa with reasonable certainty
3.	Identified by world expert in the taxa with some doubts
4.	Identified by regional expert in the taxa with high certainty
5.	Identified by regional expert in the taxa with reasonable certainty
6.	Identified by regional expert in the taxa with some doubts
7.	Identified by non-expert in the taxa with high certainty
8.	Identified by non-expert in the taxa with reasonable certainty
9.	Identified by non-expert in the taxa with some doubt
10.	Identified by the collector with high certainty
11.	Identified by the collector with reasonable certainty
12.	Identified by the collector with some doubt.

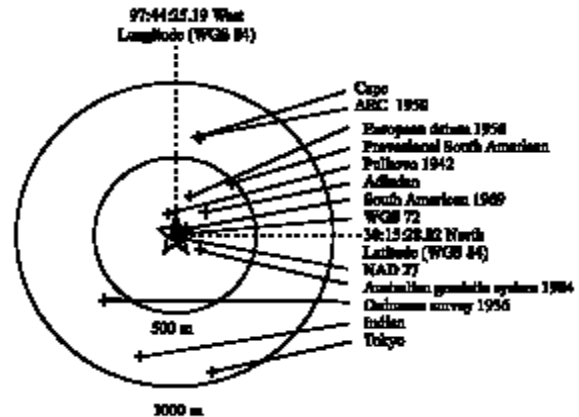


Fig. 2: Position shift from datum differences by Peter Dana (Chrisman, 1992)

A common GIS error source is that the same geographical entity (for example, a lake) may be differently described by topographers, foresters, fisheries experts, recreation specialists, wildlife officers or agronomists. This can occur on different maps, produced by different methods (remote sensing, topographic maps, aerial photographs or local surveys) and by different GIS organizations. Figure 2, for instance, presents the datum

differences for the same location. How to react to this possible lack of logical consistency in different digital data banks? Averaging routine or probability matching are common solutions.

As expected, different data dates, distinct scales and assorted class areas classification in the same layer should be avoided. Map units that are appropriate at a scale of 1:250,000, for instance, are too generalized for 1:10,000 scale map and vice-versa. The conversion of geographic coordinates (latitude and longitude) to planar ones (polar or Cartesian) by GIS input routines arise problems, as well: the mercator projection implies great distortion in high latitudes, while Azimuthal preserve neither angular nor area relationships. Yet, according to Bonham-Carter (1996), mercator is an excellent system for regions at scales of 1:250,000 and larger.

Attribute accuracy (how the estimated value approaches the true value) can be bias due to systematically wrong measurement. The international soil reference and information center showed that variation in laboratory results for the same soil samples could easily exceed $\pm 11\%$ for clay content, $\pm 20\%$ for cation exchange capacity, $\pm 10\%$ for base saturation and $\pm 0.2\%$ units for pH (Burrough and McDonnell, 1998). On the other hand, the attributes assigned to a polygon may not homogeneously to all parts of the polygon since reality is not made with black and white colors. Indeed, there is a light spectrum in the middle.

Location accuracy such as latitude, longitude or elevation above the sea level determination is a real problem, as well. Quite often, errors on map accuracy (how far is the geographical location of an object in relation to its true location on the ground) is expressed as the square root of the sum of mean deviations at different scales. Table 2 shows this relationship.

Database lineage accuracy (the ability of deriving final thematic map layers) is also a controlling concern of the growth of large spatial databases. Even, the source document is itself a distorted and abstract view of the real world is greater than digitizing errors (Painho, 1992). For the same author, in vegetation land information systems, mis-labeling of polygons (when a polygon of class α is assigned to class β) is likely to be more important than mis-location of polygon boundaries due to digitizing error. Another, traditional problem arises from keying error when the clerk entered the place code for a city as 8885 instead of 8855.

Other researchers, present spatial error in different views. According to McAlpine and Cook (1971), for instance, 2 initial maps with m_1 and m_2 segments, the overlay derived map could be estimated by $m_1 + m_2 + 2 \times \sqrt{m_1 \times m_2}$. Goodchild (1987) showed that the

Table 2: The square root of the sum of mean deviations is computed with the positional error of a set of test points, squaring the individual deviations and taking the square root of their sum

Map scale	Accepted root mean square error (m)
1:50	0.0125
1:100	0.0250
1:200	0.0500
1:500	0.1250
1:1000	0.2500
1:2000	0.5000
1:4000	1.0000
1:5000	1.2500
1:10000	2.5000
1:20000	5.0000

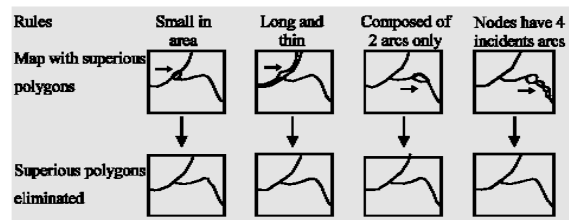


Fig. 3: General procedures to eliminate sliver polygons (Painho, 1992)

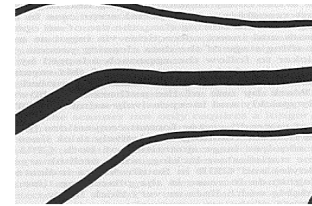


Fig. 4: Quite usual, these lines recording is left to the subjective judgment of the field worker

average number of spurious polygons may round 17%, if an overlay operation is accomplished. Chrisman (1992) discusses topological error such as dangling chain, unlabelled or conflicting labels and chain with same left/right. Confirmed by Painho (1992), in a square 100 cm map sheet size, a 0.05 mm precision represents an uncertainty of 0.05% relative to the size of the sheet. As well, major sliver polygons generation in overlay procedures may become a nightmare (Fig. 3).

A contribution from ESRI®, in particular, is the COGO module, a coordinate geometry entering procedures for land record information providing high levels of accuracy or correctness based on explicit measurement of features from some known monument but four to 20 times more expensive. Undeniably, this accuracy/cost approach may create controversial discussions. If for the overall community of municipal users the benefits seem relatively small, engineers argue that precision is necessary for survey and engineering computations.

The line width may represent a significant challenge to the vectorization designer as Fig. 4 shows Analogous,

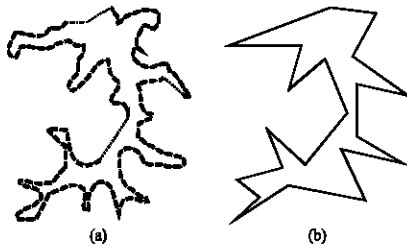


Fig. 5: The digitized lake with 204 vertices (left) versus the simplified representation by the Douglas-Peucker algorithm (right)

higher scanner resolutions are needed where lines in the map are very close together. Yet, Peuquet and Marble (1993) states that this procedure might not be necessary since, smooth lines can easily be generated by interpolation from lower precision data.

The Perkal's epsilon band model may also be useful for representation lines errors. In this paradigm, 3 situations can emerge: no error is assumed and the point is inside of the line, the observed point has a probability of 68% of lying within one standard deviation of the mean band; The point is out of the line. Still, it misses a stochastic process model regarding error accuracy. As well, the Douglas-Peucker process for line simplification suffers from the same difficulty (Fig. 5).

With the Circular Standard Error (CSE) model, a point accuracy error closely related to the true (x, y) coordinates on a topographic sheet, can be represented by a set of 2 ellipses centered on the point with one standard deviation left-right (xx direction) and above-below (yy direction). For Painho (1992), the probability that a point's true location lies somewhere within the radius circle is equal to the CSE (39.35%), that is, a 2.146 of the CSE radius will contain 90% of the distribution under the circular map accuracy standard.

The error creation by the computer word overflow was demonstrated by Gruenberger (1984). The number 1.0000001 was squared 27 times in an 80486 PC with 4 and 8 byte precision. After 27 squaring, the single precision reached a value of 8850397. With the double precision version, the final value was 674530.

According to Burrough and McDonnell (1998), the error variance in an area estimation for any polygon is given by a summation of all the errors from all the bounding cells. If m cells are intersected by the line boundary, the error variance will be

$$V = maS^4$$

where,

a = 0.0452.

V = The error variance.

S = The linear dimension of the square cell.

Other vector to raster conversion error models can be found on other literature such as Switzer (it estimates the optimal grid size and total error mismatch when converting vector polygon maps to grid ones) and Bregt (it converts data twice and compare differences).

SPATIAL INTERPOLATION UNCERTAINTY

Although, spatial data lives with uncertainty, science needs safe foundations. Uncertainty, is a dimensionless parameter for which high values are bad and lower ones are optimal. Thus, spatial uncertainty must be space geometry dependent because areas away from sample locations hold higher uncertainty. It must also take into account the variability of sample values. For obvious reasons, a particular GIS field for uncertainty computation is interpolation. Further, since, different interpolation procedures may give dissimilar results and ground truth can never be known, it may be useful to know what the predicted chance of exceeding a given upper limit is, for demand, so decisions about expensive cleanup operations can be well founded, for instance. With agricultural applications, administrators might be interested to know how much of the whole population would give a higher return than the value of a certain crop, while within environmental issues, supervisors might be looking at toxicity levels. The question is to determine how much of the population is likely to lie above or below a cut-off value.

Two issues emerge from this perspective. First, the choice of the probability threshold can be subjective. A given contamination level may be unacceptable for residential areas but tolerable for industrial yards. Second, the estimation error may be ignored. A contaminated location can be declared safe on the basis of an estimate of pollutant concentration which is incorrect but slightly less than the regulatory threshold (Goovaerts, 1997). If the true and the estimated values belong to quadrant II on an Estimated versus True grade plot, a good opportunity to invest can be missed (risk β or health cost). If it falls within quadrant IV, Fig. 6 expensive consequences can be expected (Risk α or remediation cost).

Uncertainty, about cell size and error has a close relationship with cleanup costs, too. Small blocks are desirable because cleanup can cost less money than larger ones. Nevertheless, if errors are too large then there is a realistic possibility that blocks below threshold level are treated and blocks that are really above that level are missed. According to Isaaks and Srivastava (1989), with their Walker Lake study, 7% of the area for a threshold value of a thousand was considered, while only 2% would have been suitable if the support had increased to 20x20 m².

To quantify these uncertainties becomes, then, the issue. The misclassification risk associated with a

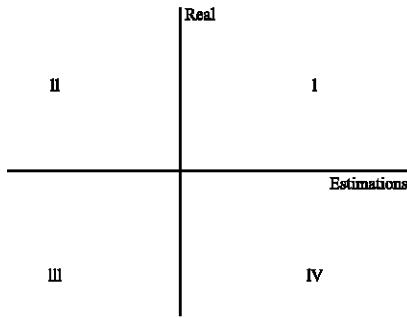


Fig. 6: False positive loss occurs when unpolluted land is classified as polluted (IV quadrant), while false negative loss when polluted land is classified as unpolluted (II quadrant)

particular physical cut-off definitely increases at threshold location boundaries since, additional sampling at those locations might be worth to avoid risk α and β costs. If the goal of a manager's decision is to minimize unnecessary cleansing and ill health costs (in conjunction with a pre-setup deterministic cost) then, it is possible to layout the total spatial health and remediation costs based on the resulting expected false negative error and false positive error models (Goovaerts, 1997).

With regard to economic land evaluation, linear programming including sensibility analysis is a possibility. Burrough (1991), presents a gradual deterministic response of pH crop illustrating soil acidity impact on crop growth: No crops, if $pH \geq 7$, Normal growth, if $pH \leq 5$ and $(7-pH)/2$ of crops, if $5 < pH < 7$. Seven years later, Burrough and McDonnell (1998), presented the Kenya annual soil erosion simulation using the:

$$\text{Universal Soil Loss} = (R_e (297 \pm 72) \times S_e (0.1 \pm 0.05) \times S_l (2.13 \pm 0.045) \times S (1.169 \pm 0.122) \times C (0.63 \pm 0.15) \times P (0.5 \pm 0.1))$$

where,

- R_e : Rainfall erosion.
- S_e : Soil erodibility.
- S_l : Symbolizes slope length.
- S : Equals slope.
- C : Cultivation.
- P : Signifies practice for rainfall-rainoff impact, yielding an annual soil loss rate of 9 ± 6 cm in 40 years.

On the basis of the assumption of null dependency among points, classical statistics presents another option:

- Estimation of the best unbiased estimator for the true population standard deviation, that is:

$$\sigma^* = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

Where,

- x_i = The *i*th sample value.
- \bar{x} = Equals the sample mean.
- n = The sample size.

- Computation of the best standard deviation confidence interval based on the true unbiased population assumption:

$$\sqrt{\frac{(n-1)\sigma^{*2}}{\chi^2_{lower}}} < \sigma < \sqrt{\frac{(n-1)\sigma^{*2}}{\chi^2_{upper}}}$$

Where ,

- σ^2 = The estimated sample variance.
- χ^2_{lower} = Equals the lowest x^2 for a certain confidence level (df = n-1).
- χ^2_{upper} = Equals the upper x^2 for a certain confidence level (df = n-1).

- Assessment of the mean confidence interval using the t-Student distribution. That is,

$$\bar{x} + t_{df,lower} \frac{\sigma^*}{\sqrt{n}} < \mu < \bar{x} + t_{df,upper} \frac{\sigma^*}{\sqrt{n}}$$

Where,

- $t_{df,lower}$ = Equals the lowest t value for a particular confidence level (df = n-1).
- $t_{df,upper}$ = Equals the upper t value for a particular confidence level (df = n-1).
- σ^* = The standard deviation dataset.

- Estimation of the probability that is likely to lie above or below a threshold value (based on the Normal distribution). This means:

$$x = \frac{\text{threshold} - \bar{x}}{\sigma^*}$$

Where,

- \bar{x} : The estimated sample mean.

The uncertainty layout of the conventional Ordinary Kriging (OK) is closely related to its variance in the following way:

$$s_e^2 = C_{00} + \sum_{i=1}^n \sum_{j=1}^n w_i w_j C_{ij} - 2 \sum_{i=1}^n w_i C_{i0} = -y_{00} - \sum_{i=1}^n \sum_{j=1}^n w_i w_j y_j + 2 \sum_{i=1}^n w_i y_{i0}$$

where,

- C_{ii} = The variance of the estimated point value.
- C_{ij} = The covariance between the *i*th and *j*th sample.
- w_i and w_j = The OK weights.
- C_{i0} = The covariance between the *i*th sample and the unknown value being estimated.

According to Soares (2000), if the sum of the weights is 1 and the average estimation error equals 0 then the Kriging error variance becomes

$$\sigma_{OK}^2 = \sum_{i=1}^n w_i \gamma(x_i, x_0) + \Psi$$

where,

- Ψ = Equals the LaGrange multiplier of the OK system.
- w_i = The OK weights while $\gamma(x_i, x_0)$ is the variance between the *i*th and the estimated point.

Hence, if errors respect the ‘bell’ curve then real values will fall within the Kriging predictor $\pm 2\sigma_{OK}$ interval for a 95% confidence level (this implies symmetry of the local distribution of errors).

However, uncertainty is not included with variogram estimation and prediction variance is underestimated. But even more critically, OK variance is not sensitive to local error for 2 major reasons: It is based on the same global variogram. Distances among locations are the only relevant factor. OK variance is mainly a geometry-dependent measure heading the assumption that an OK true error map is a better substitute. OK variance is too much of a spatial operation.

Using a theoretical dataset with 96 samples of Pb lead contamination dataset (mean = 49, variance = 457, variation coefficient = 0.43, skewness = 1.55) as an illustration purpose, Fig. 7 confirms this perspective that Kriging variance becomes higher with the absence of spatial samples and vice-versa. Notice that an anisotropic spherical model was used.

Before the final surface is produced, consistency and biased estimations absence can be tested by temporarily dropping any sample from the dataset, while its value is re-estimated using the remaining samples. Repeating this cross-validation procedure for all observations (real values versus estimated ones), it is possible to obtain the experimental mean error

$$\left(\frac{\sum (x_i^* - x_i)}{n} \right)$$

the root mean square prediction error

$$\left(\sqrt{\frac{\sum (x_i^* - x_i)^2}{n-1}} \right)$$

and the average Kriging standard error

$$\left(\sqrt{\frac{\sum \sigma_{OK}(x_i)}{n-1}} \right)$$

If the average standard error is close to the root-mean-square prediction error then the user is assessing the prediction variability correctly. However, if the average standard error is greater or less than the root

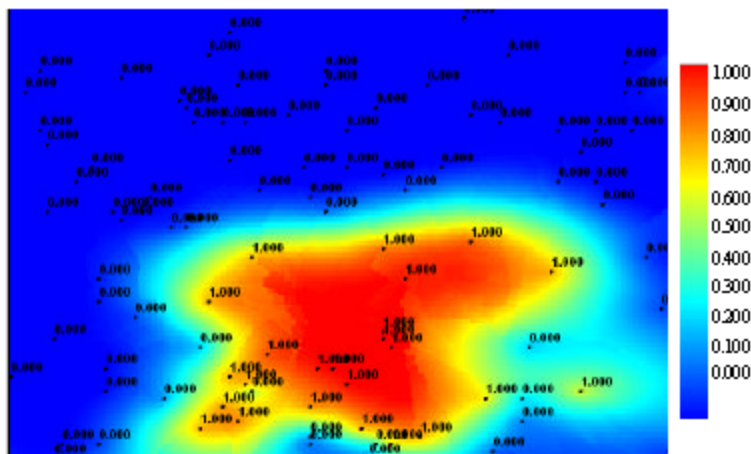


Fig. 7: The Kriging variance map

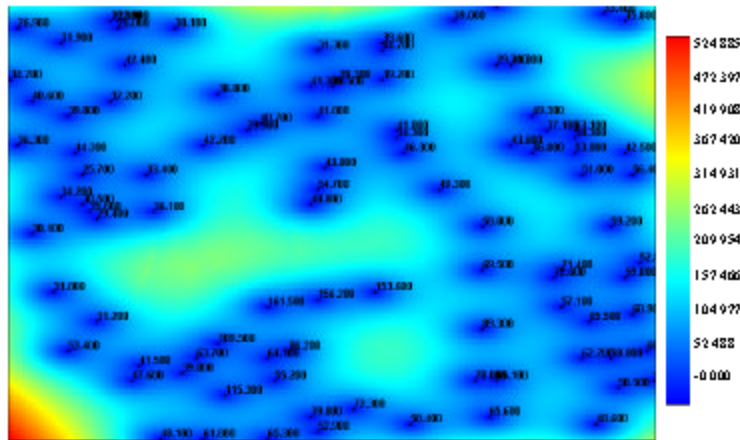


Fig. 9: Spatial estimation by Indicator Kriging (IK) for the Pb lead dataset whose cut-off level equals the 3rd quartile (59.2 ppm). As expected, the probability to exceed this particular threshold is represented by the red-orange-yellow region

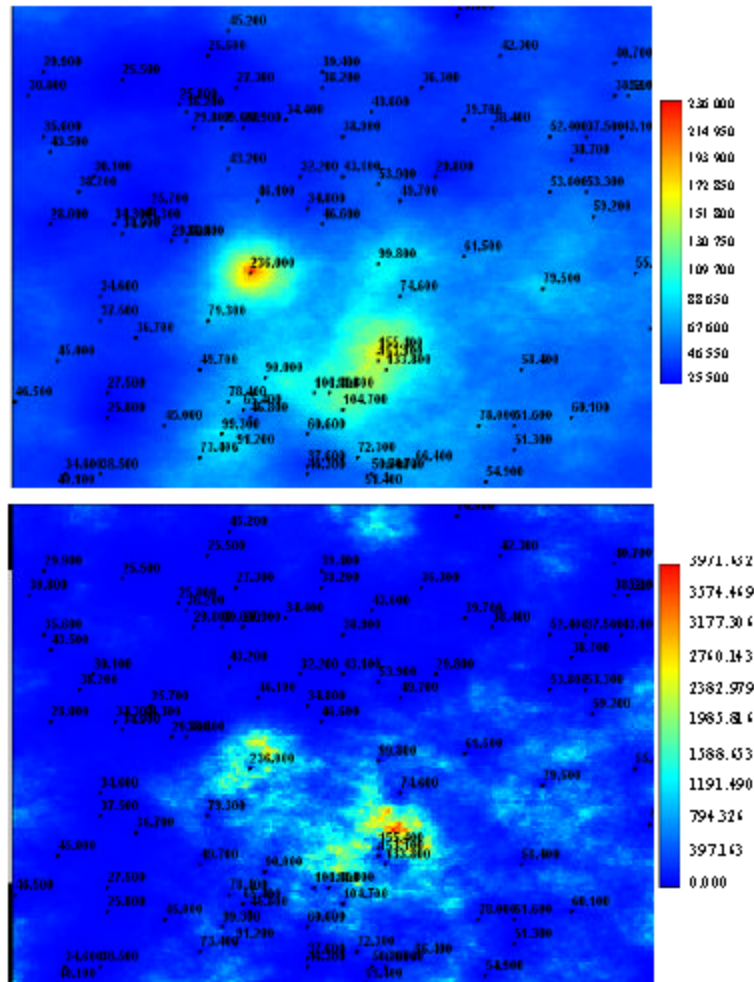


Fig. 10: The average sequential Gaussian simulation of the first 20 simulations (upper) and its variance (lower) generated by GeoMS[®] software. As expected, the simulated variance is greater where there is a major variability among the original samples or where no samples are found

configurations of possible realities, a realism issue (Kriged surface determines the most likely value at a particular location, an accuracy issue), based on the search window

Gaussian distribution but not on the optimum estimation (Wang and Zhang, 1999). For Soares (2000), the greatest potential of geostatistical simulation lies in the production of uncertainty estimations for a given cut-off value and therefore, the assessment of impact costs (Fig. 10).

CONCLUSION

If we forget the future and consider the issues that surround spatial analysis today, the gulf of ignorance between the known and the unknown is a difficulty not yet resolved. From the GIS perspective, it would be useful to have a calibrated model which could be used to describe error, to track it through GIS processes and to report uncertainty after final results are presented. From a spatial statistical view, geostatistics holds capable algorithms to handle interpolation uncertainty. Still, the faculty to deal with error outside of this field is quite lower. Geosoftware reflects this last tendency in its worst way. Researchers and commercial companies must understand that dozens of research articles are printed every year with new error evaluation methods. Regrettably, the majority of those ideas end up on the bookshelf without any practical consequence for the spatial analysis user. A theory only succeeds if it leads to a practical purpose. Categorically, the spatial error implementation is still a vacuum within most GIS products.

REFERENCES

- Bonham-Carter, G., 1996. Geographical Information Systems for Geoscientists: Modelling with GIS. 1st Edn. Pergamon Press. New York, pp: 391.
- Burrough, P. and R. McDonnell, 1998. Principles of Geographical Information Systems. 2nd Edn. Oxford University Press, pp: 333. ISBN: 0198233663.
- Burrough, P., 1991. Principles of Geographical Information Systems for Land Resources Assessment. 1st Edn. Clarendon Press, Oxford, pp: 193.
- Burrough, P., 1993. Soil variability: A late 20th century view. In: Soils and Fertilizers, 56: 529-562.
- Calkins, H., 1992. Creating Large Digital Files from Mapped Data. In Introductory Readings in GIS. In: Peuquet, D.J. and D.F. Marble (Eds.). London, Taylor and Francis, pp: 209-214.
- Chapman, A., 2008. Principles of data quality. <http://circa.gbif.net/irc/Download/kjeYAKJSmRGFqwAaUY4x8KZ1jH4pYxtv/F37w1fUI4R0AgTiySEZttf0yRVSBnGn/Data%20Quality.pdf>
- Chrisman, N., 1992. The Error Component in Spatial Data. In GIS Volume I: Principles. In: Maguire, D.J., M.F. Goodchild and D.W. Rhind (Eds.). Harlow, Longman Scientific and Technical, pp: 165-173.
- Goodchild, M., 1987. Spatial Analytical perspective on geographical information systems. In: Int. J. Geograph. Inform. Syst., 1 (4): 327-334. DOI: 10.1080/02693798708927820.
- Goovaerts, P., 1997. Geostatistics for natural resources evaluation. Oxford University Press, pp: 483.
- Gruenberger, F., 1984. Computer recreations. Sci. Am., 250 (4): 10-14.
- Isaaks, E. and R. Srivastava, 1989. In: Applied Geostatistics. 1st Edn. Oxford University Press. New York, pp: 561. ISBN: 0-19-505013-4.
- Juang, K. and D. Lee, 2000. Comparison of three Nonparametric Kriging Methods for Delineating Heavy-Metal Contaminated Soils. In: J. Environ. Quality, Madison, 29: 197-205.
- McAlpine, J. and B. Cook, 1971. Data Reliability from Map Overlay. In: Proc. 43rd Congress of the Australian and New Zealand Association for the Advancement of Science. Section 21-Geographical Sciences, Brisbane.
- Painho, M., 1992. Modelling Errors in Digital Landuse/Landcover Maps, Ph.D. Unpublished Thesis, University of California of Santa Barbara, California, pp: 140.
- Peuquet, D. and D. Marble, 1993. ARC/INFO: An Example of a Contemporary GIS. In: Readings in GIS. Peuquet, D.J. and D.F. Marble (Eds.). London, Taylor and Francis, pp: 250-285.
- Soares, A., 2000. Geostatistics for Earth Science. IST Press. 2nd Edn. Lisboa, pp: 206.
- Tomlinson, R., 1992. Current and Potential uses of GIS: The North American Experience. In: Readings in GIS, Peuquet, D.J. and D.F. Marble (Eds.). London, Taylor and Francis, pp: 203-218.
- Wang, X. and Z. Zhang, 1999. A Comparison of Conditional Simulation, Kriging and Trend Surface Analysis for Soil Heavy Metal Pollution Pattern Analysis. In: J. Environ. Sci. Health, 34: 73-89.