

Sources of Non – Stationarity in the Semivariogram

Miguel A. Cuba and Oy Leuangthong

Traditional uncertainty characterization techniques such as Simple Kriging or Sequential Gaussian Simulation rely on stationary assumptions of first and second order to describe in a numerical manner the behavior of particular variables of natural events such as a mineral deposit. But natural events are non – stationary phenomena. One of the most common solutions to this problem in practice is sub – domaining that consist of separating the mineral deposit in sub – groups that are more pseudo stationary, generally based on the physical understanding of the natural event. Even when sub – domaining has been carried on with all the available information it does not guarantee that the influence of the non – stationary features on them have been completely mitigated. In this document a way to measure the effects of non – stationarity over a domain is presented. It can be used at each stage of the sub – domaining process in order to verify how the non – stationary conditions still affects the sub - domains or to be aware of it in order to find appropriate strategies to model them. This is achieved through the interpretation of the features that the experimental semivariogram calculated via the method of moments capture from the dataset that represents a particular domain.

Introduction

The SK system relies on the semivariogram model to describe the spatial structure of the conditioning data and the semivariogram model is a function that describes the spatial continuity of a RF that obeys the intrinsic hypothesis and is first and second order stationary. The conditioning data should be part of a realization of that RF. Then the SK system will estimate the parameters of the conditional distributions under those conditions. Since SK minimizes the variance of the estimation error there is no better approach for it. Recall SGS relies on the SK system to get the realization maps.

Unfortunately the datasets sampled from natural events are non stationary, the domains are finite and the intrinsic hypothesis is vaguely satisfied. A consequence of estimating or simulating in such environment is the conditional distributions or the realization maps will be unrealistic so that they will not have any use for an economic evaluation of mineral deposit. It is necessary to minimize the impact of the lack of those conditions before proceed to apply traditional estimation or simulation techniques such as SK or SGS.

The impact of lack of the intrinsic hypothesis is small compared to the lack of the stationary conditions of the dataset and the domain. Prior of estimating or simulation the dataset should be conditioned to pseudo stationary conditions of first and second order. One way to do this is to perform sub – domaining which is to separate the domain in sub sector where the behavior of the variable to be modeled is fairly close to stationary, proceed to model it in more stable domains and join all the parts or sectors to get the final model. There also other techniques that modify the non – stationary space e.g. trend removing and transform it to a more stationary environment for modeling, and others that deal with the intrinsic assumption where they use non – parametric covariance structures.

Sub – Domaining

Uncertainty characterization of mineral deposits is an important stage in the mining industry. Geostatistical methodologies are used to build numerical models that characterize uncertainty parameters globally and locally which are required by mine planners for economic evaluation of the potential mine project. Traditional geostatistic methodologies both for estimation and simulation rely on the kriging system. The kriging system is intended to minimize the estimation variance under the assumptions of first and second order stationarity and the intrinsic hypothesis of the RF. But unfortunately the natural events that are involved in the genesis of mineral deposits are not stationary and modeling them under such assumptions could lead to wrong results.

An option to deal with this problem is to sub – divide the mineral deposit in sub – domains where the non – stationary features impact less. At this point it is important to focus on first and second order stationary features usually in Gaussian units. Ordinary kriging is widely used in mining for estimation purposes and in this case sub – domaining focus mostly on the second order stationary assumption trying to delineate

geologic unit types with high, medium and low grades of metal content taking into account the proportional effect. On the other hand if simulation is required for economic analysis first and second order stationary features in Gaussian units are focused during the sub – domaining process.

One of the common natural processes that could not be modeled under stationary assumptions is the transition of rock types such as mineralized and non – mineralized. This transition could be hard or soft. A soft transition is the gradual change in the mean and local variability of the variables being analyzed as the position changes from one rock type to another (see Figure 1 – 4); the gradual change is not necessarily linear for the mean neither for the local variability. It could also happen that the transition is abrupt which could be a consequence of structural features such as faults or post mineral geologic events. The latter one would be a scenario with presence of independent domains (see Figure 1 – 3). Most of the information that is involved in sub – domaining is the rock type logging data. From a geologic perspective rock type logging could be detailed or generalized. If detailed from a geostatistical perspective it would be possible to find domains fragmented in many rock types. On the other hand if generalized two or more geostatistical domains could be grouped in a single rock type. Even if in the generalized domain the mean grade seems to be consistent along the domain the local variability could be different and under this condition assuming stationarity would be a bad decision (see Figure 1 – 2). Figure 1 shows suitable simplistic approach for the RV behavior to model natural events and in reality they are way too much complex or intractable.

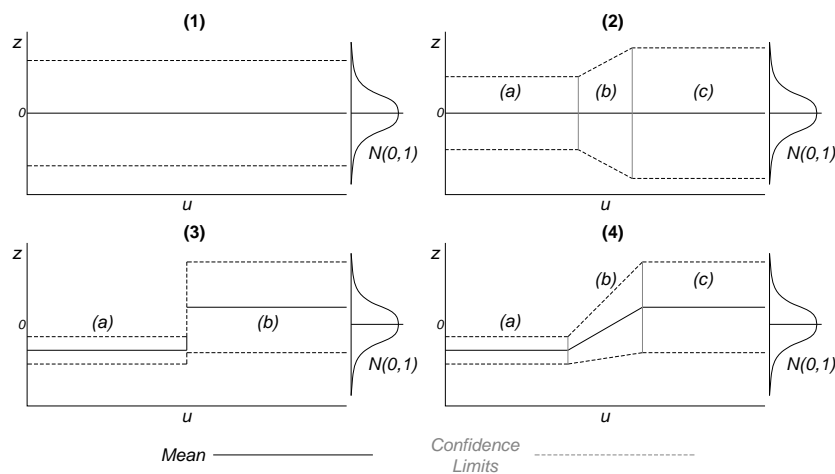


Figure 1, schematic 1D Gaussian RV environments: (1) first and second order stationary; (2) first order stationary but locally variable variance with two first and second order stationary domains (a), (c) and one first order stationary with variable local variance (b); (3) non – stationary with two independent first and second order stationary domains (a) and (b); (4) non – stationary with two first and second order stationary domains (a) and (c) and one non – stationary domain (b).

In estimation or simulation the aim is to assess uncertainty, it is achieved through the calculation of the mean and variance at the locations of unknown values conditioned to available data from the domain. Estimating or simulating in a non – stationary environment does not guarantee that the obtained results are correct for modeling natural events. Since SK and consequently SGS rely on stationary and intrinsic hypothesis assumptions, even when the estimated or simulated mean values could be somehow correct due to strong presence of conditioning data the estimation variance is not correct, it could be underestimated or overestimated according to the underlying behavior of the variables in the domain (see Figure 1). On the other hand the intrinsic hypothesis impact could be considered small compared with the impact due to the lack of pseudo stationary conditions.

Semivariogram in a Stationary and Non – Stationary Environments

The variogram is the measure of variability between two point values $z(\mathbf{u}_1)$ and $z(\mathbf{u}_2)$ at two different locations. (Barry 2004) Analytically the variogram represent the variance structure between two RV at different locations. It can be expressed as (Cressie, 1993) the variance of the difference of two RV spatially located at locations \mathbf{u}_1 and \mathbf{u}_2 (1). In order to calculate it many realizations of the two random variables

$Z(\mathbf{u}_1)$, $Z(\mathbf{u}_2)$ would be required, unfortunately in practice only one is available which is the conditioning data.

$$2\gamma(\mathbf{u}_1, \mathbf{u}_2) = \text{var}\{Z(\mathbf{u}_1) - Z(\mathbf{u}_2)\} \quad \forall \mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^d \quad (1)$$

To overcome that problem the intrinsic hypothesis is assumed. (Matheron 1971) The RF Z obeys the intrinsic hypothesis if the variance of the spatial difference $Z(\mathbf{u}) - Z(\mathbf{u} + \mathbf{h})$ exist and is independent of location \mathbf{u} but are function of the separation vector \mathbf{h} . The variance of the spatial difference is the variogram expression (2) and the expected value is a function of the separation vector and assumed to be a linear drift (3). Removing the linear drift from $Z(\mathbf{u})$ for all the locations in the form $R(\mathbf{u}) = Z(\mathbf{u}) - m(\mathbf{u})$ the expected value of the spatial difference of the residuals $R(\mathbf{u})$ become zero and the variogram expression can be written as the expected value of the squared spatial difference (4).

$$\text{var}\{Z(\mathbf{u}) - Z(\mathbf{u} + \mathbf{h})\} = 2\gamma(\mathbf{h}) \quad (2)$$

$$E\{Z(\mathbf{u}) - Z(\mathbf{u} + \mathbf{h})\} = m(\mathbf{h}) \quad (3)$$

Removing the drift component from the RV $Z(\mathbf{u})$ is equivalent to make the RF first order stationary. Recall the RF R should obey the intrinsic hypothesis.

$$R(\mathbf{u}) = Z(\mathbf{u}) - m(\mathbf{u})$$

Then the variogram is expressed as follows:

$$2\gamma(\mathbf{h}) = E\left\{\left[R(\mathbf{u}) - R(\mathbf{u} + \mathbf{h})\right]^2\right\} \quad (4)$$

The difference between expression (2) and (4) is in the mean value of the RV. (Gneiting and others, 2001) regardless of the definition of variogram or semivariogram due to the 0.5 factor, they distinguish between centered and non – centered variograms. Expression (2) is called centered and (4) is non – centered. In order to link the semivariogram expression with the covariance and variance it is necessary to make assumptions of second order stationarity in $R(\mathbf{u})$. Then the semivariogram can be written in terms of the difference between the variance and the covariance (5).

$$\gamma(\mathbf{h}) = \sigma^2 - c(\mathbf{h}) \quad (5)$$

$$\gamma(\mathbf{h}) = \sigma^2 (1 - \rho(\mathbf{h})) \quad (6)$$

In equation (5) $\gamma(\mathbf{h})$ is the semivariogram, σ^2 is the variance and $c(\mathbf{h})$ the spatial covariance and in expression (6) $\rho(\mathbf{h})$ is the spatial correlation coefficient. The semivariogram, the spatial covariance and the spatial correlation coefficient are function of the separation vector \mathbf{h} . Recall the spatial covariance at $\mathbf{h}=0$ is the variance therefore at no spatial correlation the semivariogram takes the value of the variance. The spatial correlation coefficient at $\mathbf{h}=0$ is 1.

Due to the Intrinsic Hypothesis the kriging estimation variance is dependent of the spatial configuration of the conditioning data rather than variability patterns as a function of location \mathbf{u} in the domain. In a mining project analysis the estimation variance is often used as a parameter to rank uncertainty so this value could be misleading if not used properly.

Since first and second order stationary assumptions and intrinsic hypothesis are assumed these three features have an impact on the modeling of the required variables and uncertainty characterization in the domain. The natural event has to be assumed that follows a behavior like Random Function under the three previous assumptions. Unfortunately the spatial features of the natural events do not follow any of those assumptions and they have to be accommodated in a suitable manner so that the non – stationary features do not affect too much the domain. Prior to modeling it is important to make the big domain as much stationary as possible, some techniques to achieve this are sub – domaining and trend modeling. And to deal with non – stationary covariance or semivariogram functions many approaches have been developed, see for example Sampson and Guttorp 1992, proposed to move from an initial “geographical” dimension to a “dispersion” dimension where the intrinsic hypothesis is more suitable.

In practice expression (1) makes more sense in order to model natural events. The spatial variability is a function of the location and between two locations. One simple example is the presence of anisotropic behaviors where the spatial correlation in the direction of major continuity is different to the minor's. The use of anisotropic ratios is an attempt to convert from a space where the intrinsic assumption is not valid to another one where it is. But natural events are more complicated than this. In addition to that there is no condition in the natural events that could explain that the anisotropic behavior is elliptical. It is just a simplistic approach to average local spatial continuity structures.

Impact of non – stationarity in experimental semivariogram calculation

From expression (4) the estimator of the semivariogram is (7) also known as method-of-moments (Matheron 1962) which can be also interpreted as the average value of the orthogonal distances of the data pairs to the 45° line in the h – scatter plot.

$$\hat{\gamma}(\mathbf{h}) = \frac{1}{2n(\mathbf{h})} \sum_{i=1}^{n(\mathbf{h})} [r(\mathbf{u}_i) - r(\mathbf{u}_i + \mathbf{h})]^2 \quad (7)$$

Where $n(\mathbf{h})$ is the number of available pairs for the separation vector \mathbf{h} . $r(\mathbf{u}_i)$ and $r(\mathbf{u}_i + \mathbf{h})$ are the i -th data pair of the initial dataset at the head and tail of the separation vector \mathbf{h} respectively.

The experimental semivariogram plot consists of many different experimental semivariogram values calculated for different separation vector ($\mathbf{h}_i, i=1, \dots, n$) in ascending order of distance for one single direction for directional semivariogram plots and in all direction for omnidirectional semivariogram plots. In presence of sparse data the separation vector \mathbf{h} is took with tolerances. In order to calculate the experimental semivariogram the initial dataset r_0 with mean m_0 and standard deviation s_0 is split in two parts or sub datasets $r_{\mathbf{u}}$ and $r_{\mathbf{u}+\mathbf{h}}$ with mean $m_{\mathbf{u}}$, $m_{\mathbf{u}+\mathbf{h}}$ and standard deviation $\sigma_{\mathbf{u}}$, $\sigma_{\mathbf{u}+\mathbf{h}}$ respectively which correspond at the two extremes of the separation vector \mathbf{h} . The representative dataset $r_{\mathbf{h}}$ for each separation vector \mathbf{h} consist of the union of the two sub datasets $r_{\mathbf{h}} = r_{\mathbf{u}} \cup r_{\mathbf{u}+\mathbf{h}}$ and for small separation distances the initial dataset and the representative dataset tend to be the same $r_0 \approx r_{\mathbf{h}}$ but as the distance of the separation vector increases they tend to be different and $r_0 \neq r_{\mathbf{h}}$ since the sub datasets have less information and the experimental semivariogram value is less reliable. Notice the initial dataset r_0 is not influenced by any declustering weights, it is assumed to be fairly representative of the domain.

At this point the available dataset for calculating the experimental semivariogram is limited or finite and the number of available pairs of samples for each configuration of the separation vector is an issue now. Even though stationarity is assumed in a finite domain the mean and variance are not constant for the entire domain they are a function of how representative the sub datasets at the head and tail of the separation vector \mathbf{h} are. These are some features of finite domains that cause problems with the stationary semivariogram features, like the number of available samples for each separation vector and presence of sparse sampling data.

From Figure 2 when Gaussian space is assumed the calculation of the semivariogram estimator for each separation vector \mathbf{h} to be according to expression (7) should produce a symmetric h – scatter plot and the data pairs should follow a bivariate Gaussian distribution. In presence of enough data the two sub datasets are expected to have similar, means equal to zero, variances and marginal distribution shapes closely to Gaussian shape, equal to one. In general the only parameter that is expected to change is the spatial correlation coefficient $\rho(\mathbf{h})$. In this context the notice that the covariance plot is similar to the correlograms. Unfortunately those required conditions are not present in real data but their impact can be quantified numerically through the semivariogram expression.

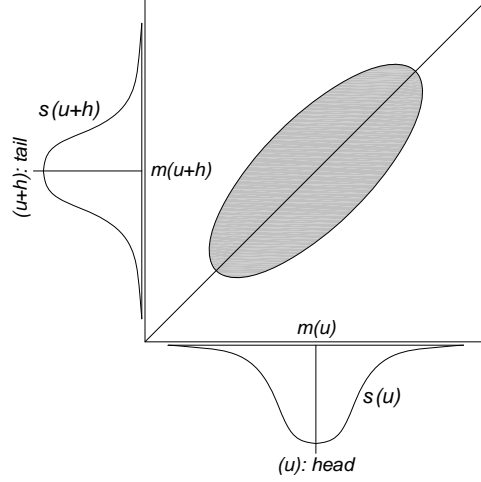


Figure 2, ideal h – scatter plot of semivariogram calculation for a dataset which follows a Gaussian distribution. In Gaussian units the data pairs should lie in a Gaussian bivariate joint distribution and the spatial continuity is measured by the correlation coefficient which is a function of the separation vector. m_u , m_{u+h} are the means and σ_u , σ_{u+h} are the standard deviations of the distributions of the sub – datasets at the extremes of the separation vector h .

The influences of the mean and variances of the sub datasets at the two extremes of the separation vector \mathbf{h} can be calculated from expression (7). Expanding the initial expression then adding and subtracting their respective average values of the sub datasets r_u and r_{u+h} it is possible to express (7) in terms of covariance, means and variances of the two sub datasets as follows:

$$\begin{aligned}\hat{\gamma}(\mathbf{h}) &= \frac{1}{2n(\mathbf{h})} \sum_{i=1}^{n(\mathbf{h})} \left[(r(\mathbf{u}_i))^2 - 2r(\mathbf{u}_i)r(\mathbf{u}_i + \mathbf{h}) + (r(\mathbf{u}_i + \mathbf{h}))^2 \right] \\ &= \frac{1}{2n(\mathbf{h})} \sum_{i=1}^{n(\mathbf{h})} \left[(r(\mathbf{u}_i))^2 + (r(\mathbf{u}_i + \mathbf{h}))^2 \right] - \frac{1}{n(\mathbf{h})} \sum_{i=1}^{n(\mathbf{h})} [r(\mathbf{u}_i)r(\mathbf{u}_i + \mathbf{h})] \\ &= \frac{1}{2} [\sigma_u^2 + \sigma_{u+h}^2] - c_h + \frac{1}{2} [m_u - m_{u+h}]^2\end{aligned}$$

The semivariogram expression is written as a function of the differences of standard deviations and difference of means of the sub datasets at locations \mathbf{u} and $\mathbf{u}+\mathbf{h}$ and the covariance (8).

$$\hat{\gamma}(\mathbf{h}) = \frac{1}{2} [\sigma_u - \sigma_{u+h}]^2 + \frac{1}{2} [m_u - m_{u+h}]^2 + \sigma_u \sigma_{u+h} - c_h \quad (8)$$

Or the correlation coefficient (9) for a separation vector \mathbf{h} .

$$\hat{\gamma}(\mathbf{h}) = \frac{1}{2} [\sigma_u - \sigma_{u+h}]^2 + \frac{1}{2} [m_u - m_{u+h}]^2 + \sigma_u \sigma_{u+h} (1 - \rho_h) \quad (9)$$

From equation (8) the semivariogram expression consists of four components. (a) Half of the average of the squared difference of their respective standard deviations plus (b) half of the average of the squared difference of their respective means plus (c) the product of the two standard deviations of the sub datasets of the head and tail of the separation vector minus (d) the covariance of the two sub datasets separated by the \mathbf{h} vector. Notice that when the means and variances of the distributions of the two datasets are fairly equal ($\sigma_u \approx \sigma_{u+h} \approx \sigma$ and $m_u \approx m_{u+h} \approx m$) the experimental semivariogram expression (8) relies on the spatial covariance and in expression (9) on the spatial correlation coefficient which are the stationary forms of the semivariogram.

In a stationary environment negative correlation coefficients are interpreted as if the initial dataset is affected by a large trend pattern ($m_u \neq m_{u+h}$) and are thought to be the reason why the experimental semivariogram values take values greater than the variance of the dataset. (Gringarten & Deutsch 2001) Trends in the data can be identified from experimental semivariogram, which keeps increasing above the theoretical sill. In simple terms, this means that as distances between data pairs increase the differences between data values systematically increase. But in practice this is not necessarily correct since the correlation coefficient $\rho(\mathbf{h})$ is a function of the two sub – datasets separated by \mathbf{h} , the impact occur when differences in mean and variances of the two datasets are different (9).

In presence of sparse data the data pairs are built using tolerances in the search of pairs, this makes that for one data point for a particular lag distance it could be paired with more than one sample. In that case the distribution of the sub dataset at the head of the separation vector \mathbf{h} is the result of the weighting each sample location by the number of repetitions in the pairing due to the search tolerances. This value is not fairly correct, it depends on the tolerances. For small tolerances the semivariogram calculation is more representative. Big tolerances tend to mask spatial features of the dataset e.g. hiding of the anisotropic features of the domain.

Case Study 1

The case study is a 1D unconditional simulated dataset of 1000 data points regularly spaced with a spherical semivariogram model of 25 units of range. Three different trend cases were added to the initial dataset in order to describe them via the experimental semivariogram expression. The first trend component is a linear trend in the form $y=ax+b$ with a negative slope along all the dataset. The second trend component is a symmetrical convex shape trend. For the third case the mean component remains the same but the variance of the two halves of the dataset were modified (see Figure 3).

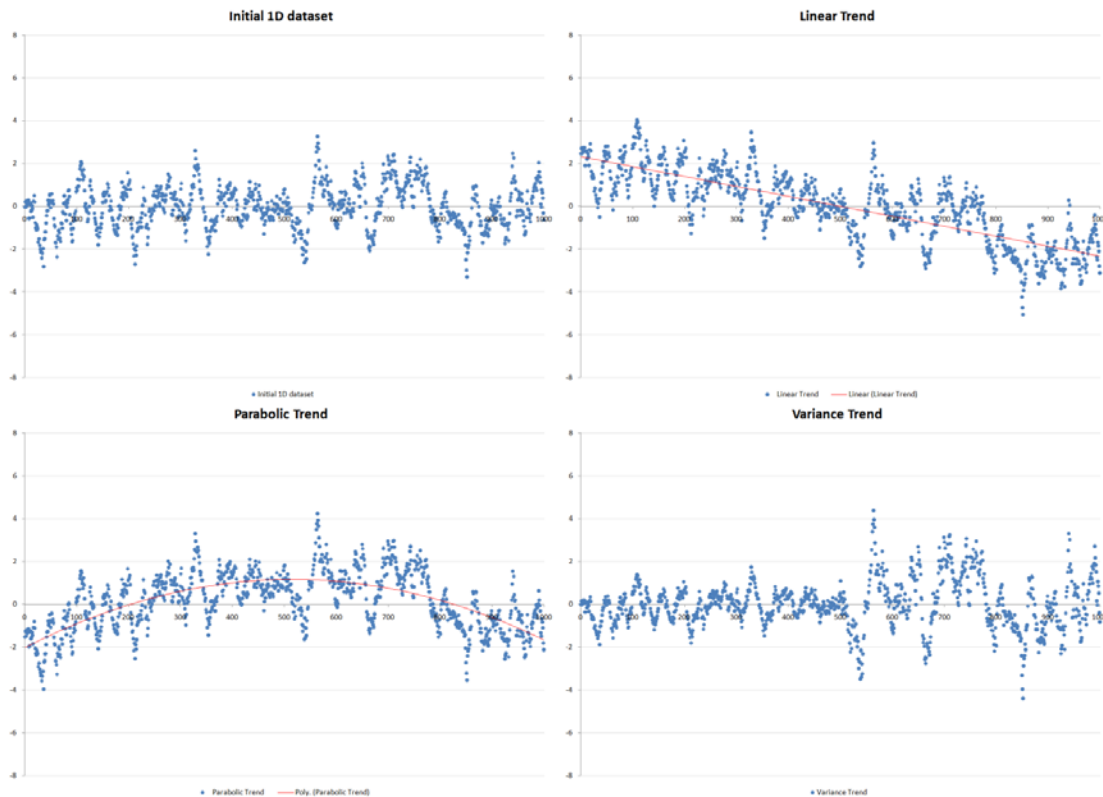


Figure 3, Initial dataset influenced by linear trend, parabolic trend and local variability in variances.

Experimental semivariogram values can be calculated using expressions (8) or (9) so that the semivariogram value is divided in three parts: 1) the mean component, 2) the variance component and 3) the stationary component. The semivariogram plots for each case are calculated for the half of the size of

the domain as a maximum separation distance in order to have enough data pairs that represent the dataset. The contributions of each component of the experimental semivariogram are represented as regions of different color (see Figure 5).

Even when the initial dataset is supposed not to have any influence of mean trend or variance trend it can be seen that there is presence of those components in the semivariogram plot. This effect can be seen comparing the two distributions of the first and second half of the entire initial dataset (see Figure 4). Notice that even when the initial dataset lies on standard normal distribution $N(0,1)$ it does not necessarily mean that locally is the same and that behavior can be captured by the experimental semivariogram. For comparison those variation in the initial dataset are assumed negligible, but that will be captured by the experimental semivariogram.

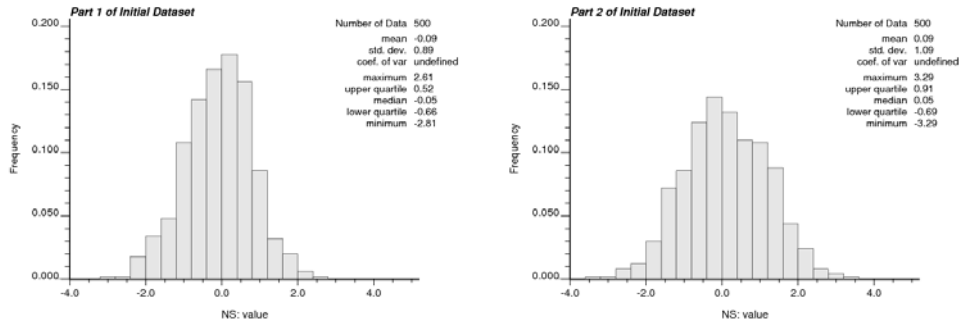


Figure 4, distribution of the first 500 data points in the dataset (left) and 500 last data points (right)

In three out of the four cases the experimental semivariogram capture the first and second order non – stationarity effect in their respective datasets. There is particular case where the contribution of the mean trend component cannot be seen in the experimental plots. It happens when the mean trend component is symmetric (the convex shape trend). The mean difference component occur when differences of the means appear as the separation distance increases and in this case since the trend component is symmetric the means of the sub datasets are cancelled ($m_u = m_{u+h}$). On the other hand notice that under that condition the only component that remains is the stationary component and the trend is capture under the stationary conditions of semivariogram (experimental semivariogram values greater than the variance of the initial dataset). This condition is interesting since in presence of directional symmetric mean trend shapes the stationary conditions of the semivariogram are valid again. This can also be seen in the h – scatter plots for each case (see Figure 8, Figure 9, and Figure 10)

Case Study 2

The second case study was run with real data for a vertical direction. The Amoco3d.dat dataset consist of 62 vertical wells with on average 53 samples per well (35 minimum and 66 maximum). There are a total of 3303 available data points in the dataset. Vertically the separation distance between the samples is one unit of distance.

A vertical semivariogram plot of the porosity variable in normal score units is calculated (Figure 6). Since the dataset is regularly spaced in the vertical direction a lag distance interval of one unit of distance is chosen so there will not be overlap between pairs. Traditionally it could be inferred that there is a presence of a mean trend since the vertical variogram does not reach the sill until the lag interval 38 where $\gamma(38) = 0.9280$ with 950 pairs out of 3241 but there are two trend components that are captured for the dataset.

The two trend components are present in the initial dataset for the vertical direction (variance trend and mean trend) (see Figure 7) and the variance for each location of the semivariogram start to be different from the lag 25. It is also important to verify how representative the sub datasets are for each lag interval.

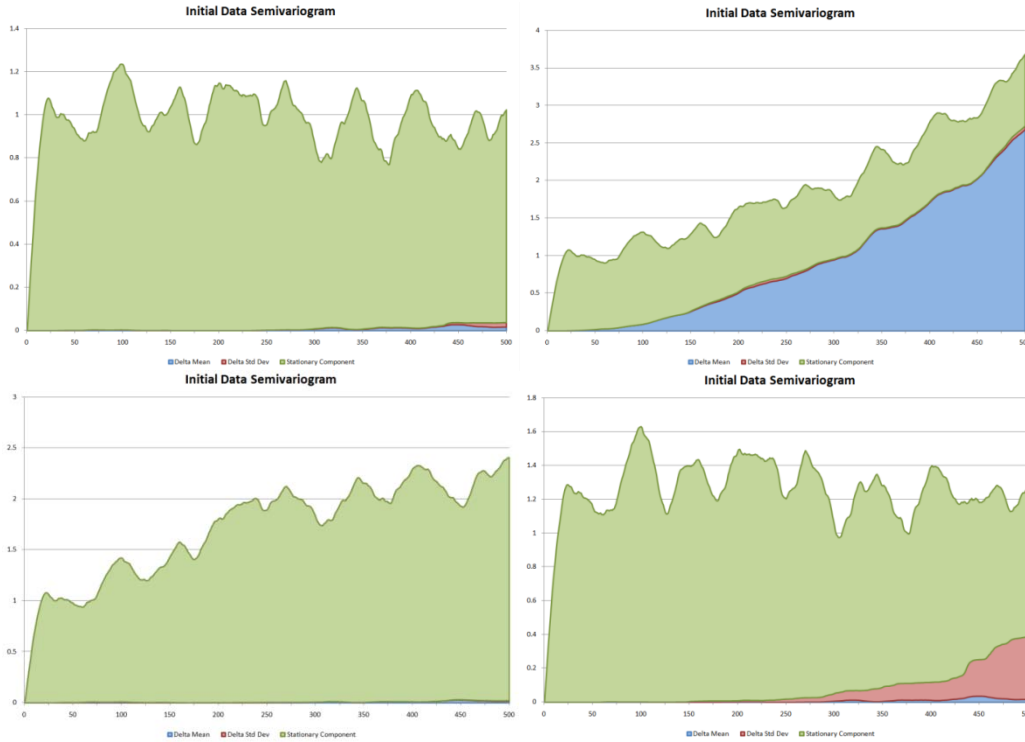


Figure 5, experimental variograms for initial dataset and influenced by linear trend, parabolic trend and locally variable variance. In the four plots the blue region represents the variation of the mean component, the red region the variation of the variance component and the green region the stationary component of the experimental semivariogram value.

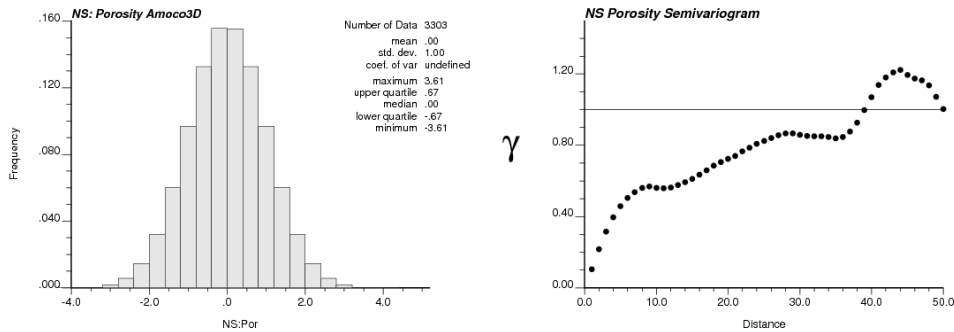


Figure 6, distribution (left) and vertical experimental variogram (right) for the NS values of porosity

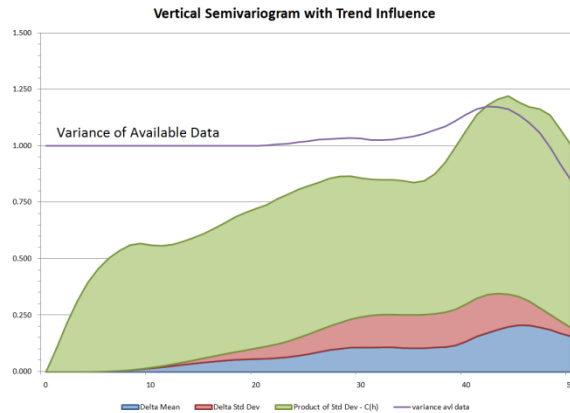


Figure 7, experimental semivariogram of vertical direction for oilsand dataset

Conclusions

The SK system is unbeatable calculating the parameters of the conditional distribution at a location where the variable is unknown conditioned to previous existing data. Traditionally the SK system uses a spatial variability model that describes the behavior of a RF which is first and second order stationary and obeys the intrinsic hypothesis. The estimated parameters of the conditional distribution then are also constrained to those conditions; additionally the conditioning dataset should be part of the previously mentioned RF. From those conditions the stationary assumptions are usually more important than the intrinsic hypothesis assumption since in presence of strong conditioning data the impact of the latter can be considered as negligible.

In a geostatistical analysis the aim is to assess uncertainty. Usually parameters such as mean and variance are enough to describe it. In some cases when only the mean is calculated or considered important it could lead to a misinterpretation of the geostatistical model results and consequently lead to a wrong decision making process in a economic evaluation of a potential mineral deposit project. In that case the geostatistical results are obtained but not interpreted correctly.

Since traditional estimation on simulation techniques rely on stationary assumptions and sampled datasets from natural events are not necessarily stationary this approach is a useful tool to measure numerically the effects of non – stationarity via experimental semivariogram calculation in the datasets and proceed to make decisions about sub domaining or any other data processing techniques in order to make the estimation or simulation more consistent with the natural event being modeled.

Non stationary characteristics present in real datasets that represent natural events can be captured from it through interpretation of the many different features of the experimental semivariogram tool comparing how different is the calculated from the ideal case. The semivariogram contains enough information to describe the spatial continuity beyond the stationary assumptions.

References

- Chilès J.P. & Delfiner, P. *Geostatistics: modeling spatial uncertainty*. Wiley-Interscience, New York, 1999.
- Genton M.G. *Variogram Fitting for Generalized Least Squares Using an Explicit Formula for the Covariance Structure*, *Mathematical Geology*, Vol. 30, No. 4, 1998.
- Goovaerts, P. *Geostatistics for Natural Resources Evaluation*. Oxford University Press, New York, 1997.
- Matheron G. *The Theory of regionalized Variables and Its Applications*, Ecole Nationale Supérieure de Paris, Paris 1971
- Journel A.G. & Huijbregts Ch.J. *Mining Geostatistics*. The Blackburn Press, 1978.
- Gringarten & Deutsch, *Variogram Interpretation and Modeling*, *Mathematical Geology*, Vol. 33, No. 4, 2001.

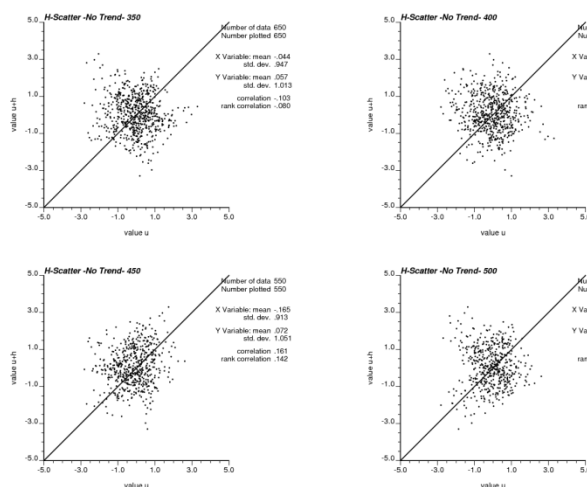


Figure 8, h – scatter plots for the initial case with no influence of trend.

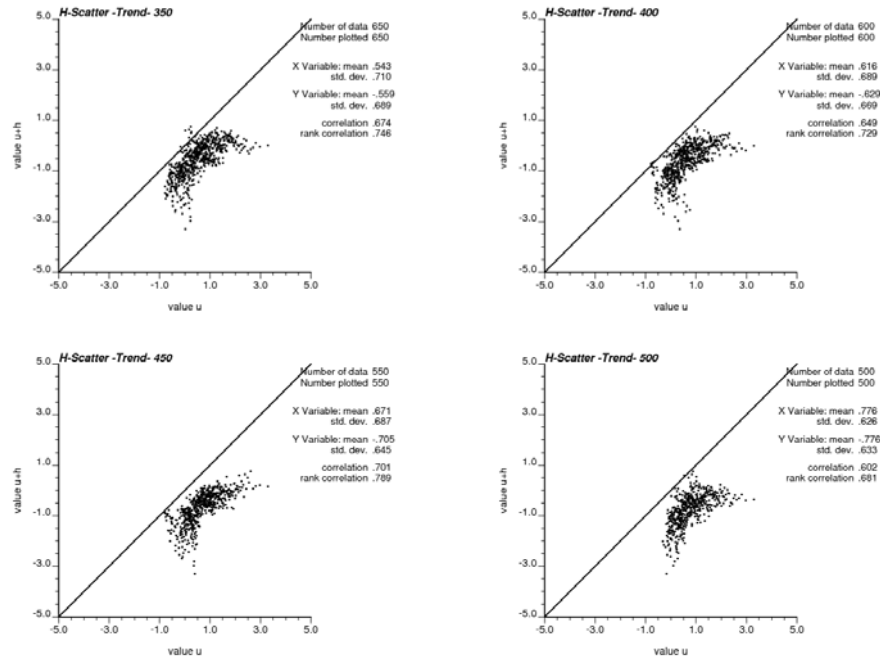


Figure 9, h – scatter plots for the linear trend case

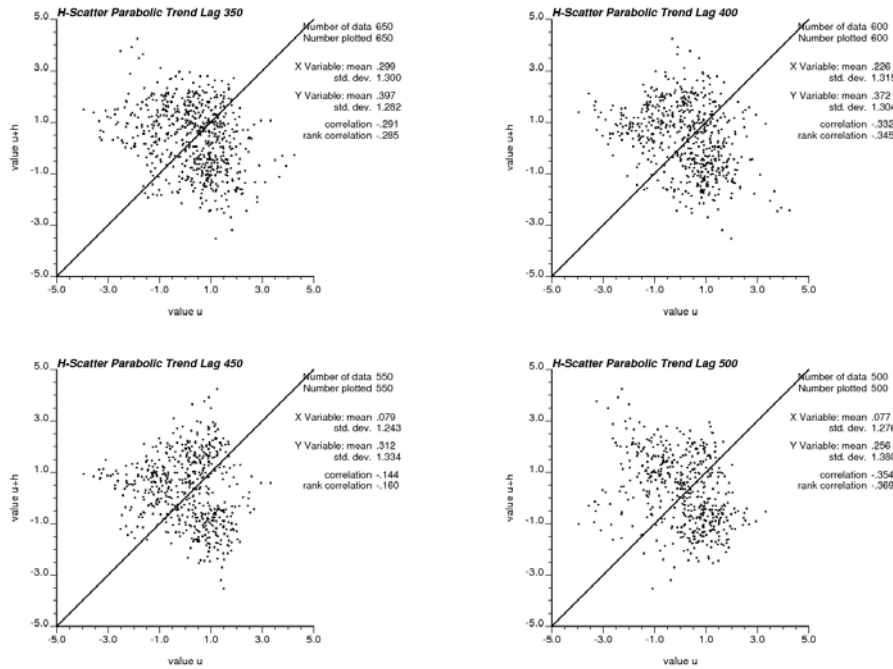


Figure 10, h – scatter plot for the parabolic shape trend case