A Generalized Direct Lognormal Simulation Algorithm with Software: slogsim

John Manchuk, Oy Leuangthong and Clayton V. Deutsch

Centre for Computational Geostatistics Department of Civil & Environmental Engineering University of Alberta

The idea of direct simulation is becoming more established with the incorporation of unstructured grids in ore body and reservoir modeling. Direct kriging and simulation permits reliable integration of multiscale data. Simply performing kriging on data in original units, however, leads to a variance that is incorrect as real data exhibit a proportional effect. This paper introduces a direct simulation algorithm for data that appears lognormally distributed. Two types of direct simulation can be executed: (1) a naïve type that performs kriging and simulation on the data directly and (2) a form that corrects the kriging variance according to the proportional effect inherent in lognormal data. In fact, the proposed algorithm is not direct simulation – it is not based on the simple kriging (SK) principle that underlies the theory and publications of direct simulation. The proposed algorithm considers links the kriging estimate and the kriging variance according to the lognormal model; the SK principle requires independence of the estimate and variance. Software is developed and documented.

Introduction

Direct sequential simulation (DSS) [1, 2, and 3] has been proposed because of its ability to account for data of various support volumes and populate unstructured grids. Kriging and simulating in original units is the essential idea of DSS. Although kriging provides a valid estimate and variance for a conditional distribution, the resulting homoscedastic variance poses a significant problem when original data units are considered; the uncertainty in low-valued areas is over stated and the uncertainty in high-valued areas is understated.

Real data often exhibit a classical heteroscedastic relationship between the local mean and variance, commonly referred to as the proportional effect [4]. With kriging as the main engine in DSS, the resulting simulated values do not reproduce a heteroscedastic feature; a method must be developed to account for the proportional effect inherent in original data units.

Simple kriging (SK) is important in DSS because of its ability to reproduce the covariance even if the conditional distributions are not Gaussian [1]. Covariance reproduction using SK can be easily demonstrated; however, it only holds if the variance of the data is homoscedastic. In the case of lognormal data the variance is heteroscedastic.

This paper proposes a solution to the homoscedastic kriging variance problem of DSS by introducing a direct lognormal simulation algorithm. This is a particularly interesting case since the mathematical relationship between the lognormal and the commonly used normal distribution is well known, as are the equations that describe the proportional effect of lognormal data [5]. Knowing these relations, the kriging variance can be calibrated to honor the heteroscedasticity inherent in lognormal data. This well posed case provides valuable insight into the nature of DSS.

Theoretical Background

A key requirement in the derivation of DSS is that the variance be independent of the conditional mean; however, the lognormal model is unique in that this requirement is *not* met; it is a different model. Working with a lognormal distribution was chosen because there is an analytical relationship between Normal-(0,1) data as is used in Sequential Gaussian Simulation, and lognormal data. Also, many actual data sets are indicative of lognormal distributions. Analytical transformations can be carried out with the data distributions as well as with the variograms.

The definition of a lognormally distributed variable is as follows: A random variable, $Z \mid z(\mathbf{u}) > 0$, is lognormal with a mean *m* and standard deviation σ if the natural logarithm of $Z(\mathbf{u})$ is normally distributed with mean α and standard deviation β . Knowing the relation between $Z(\mathbf{u}) \rightarrow logN(m,\sigma)$ and $X(\mathbf{u}) \rightarrow N(\alpha,\beta)$ one can transform between Gaussian and lognormal distributions. Equations 1, 2 and 3 show the relationship between $X(\mathbf{u})$, $Y(\mathbf{u})$, and $Z(\mathbf{u})$, where $Y(\mathbf{u})$ is a standard normal distribution. Equations 4 and 5 show the relationship between *m* and σ of $Z(\mathbf{u})$ with α and β of $X(\mathbf{u})$.

$$X(\mathbf{u}) = \alpha + \beta \cdot Y(\mathbf{u}) \tag{1}$$

$$Z(\mathbf{u}) = e^{X(\mathbf{u})} \tag{2}$$

$$Z(\mathbf{u}) = e^{\alpha + \beta \cdot Y(\mathbf{u})}$$
(3)

$$\alpha = \ln(m) - \frac{\beta^2}{2} \tag{4}$$

$$\beta^2 = \ln\left(1 + \frac{\sigma^2}{m^2}\right) \tag{5}$$

Equations 6 and 7 describe the normal and lognormal probability distribution curves. Figure 1 shows the change in the distribution shapes as $Y(\mathbf{u})$ is converted into $X(\mathbf{u})$ and as $X(\mathbf{u})$ is transformed into $Z(\mathbf{u})$.

$$f(x) = \frac{\exp\left[-\frac{1}{2}\left(\frac{\ln(x) - \alpha}{\beta}\right)^2\right]}{\beta \cdot x\sqrt{2\pi}}$$
(6)

$$g(x) = \frac{\exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]}{\sigma\sqrt{2\pi}}$$
(7)

As mentioned, there is an analytical relationship between the variogram of a Gaussian variable and a lognormal variable. If the variogram in Gaussian space is known, it can be converted to the variogram in lognormal space through the use of Equation 8. Figure 2 shows a spherical variogram for Gaussian data, the corresponding lognormal variogram and the difference between them.

$$\gamma_{Z}(h) = 1 - \frac{m^{2}}{\sigma^{2}} \left[e^{\beta^{2} \cdot (1 - \gamma_{Y}(h))} - 1 \right]$$
(8)

Where $\gamma_Z(h)$ is the variogram of the lognormal variable and $\gamma_Y(h)$ is the variogram of the Gaussian variable.

An additional characteristic of lognormal distributions that we must account for is the proportional effect. Before moving onto the equation, some comments on the conditional distributions resulting from kriging must be made. In simulation, our goal is to calculate a conditional distribution (conditional to some number of local data) for simulation. In a multivariate Gaussian case we transform the data, infer the parameters in Gaussian units, and then back transform the result. The transform and back transform are particularly easy when the data are lognormal. In fact, the shapes of all conditional distributions in original units are lognormal when the original global histogram is lognormal (see Appendix). The key idea of DSS is to krige in original units, but we must establish the correct variance, which is heteroscedastic, that is it depends on the magnitude of the data and estimate. The heteroscedasticity or proportional effect is automatically accounted for in the back transform. As there is no back transform in DSS we have to build in some form of correction. What makes the lognormal case unique is that we know the proportional effect analytically and it is derived from Equation 5:

$$\beta_L^2 = \ln\left(1 + \frac{\sigma_Z^2}{m_Z^2}\right)$$

$$e^{\beta_L^2} = 1 + \frac{\sigma_Z^2}{m_Z^2}$$

$$\sigma_Z^2 = m_Z^2 (e^{\beta_L^2} - 1)$$
(9)

Where β_L^2 is the homoscedastic kriging variance in $N(\alpha,\beta)$ units and m_Z and σ_Z^2 are the estimate and variance from kriging data in original units.

Since m_z is the estimate from kriging it can be denoted by $z^*(\mathbf{u})$. From Equation A2 in the Appendix it is shown that:

$$\beta_L^2 = \beta_G^2 \cdot \sigma_Y^2$$

where β_G^2 is the global variance of $X(\mathbf{u})$ and σ_Y^2 is the Gaussian kriging variance. Substituting these results into Equation 9 yields Equation 10:

$$\sigma_{Z,C}^{2} = \left[z^{*}(\mathbf{u})\right]^{2} \left(e^{\beta_{G}^{2} \cdot \sigma_{Y}^{2}} - 1\right)$$
(10)

Where $\sigma_{Z,C}^2$ is the corrected variance, σ_Y^2 is the local variance in Gaussian space, and β_G^2 is the global variance of ln(Z). A major implication from Equation 10 is that kriging would have to be performed twice; once to get the kriging variance in Gaussian space (σ_Y^2) and again to get the estimate in lognormal space ($z^*(\mathbf{u})$).

To experimentally show this relation, a Gaussian variable was generated using unconditional simulation and the corresponding lognormal values were calculated via Equation 3 ($m=\sigma=100$).

Kriging was performed on both data sets and the GAM program [6] was run on the results at a lag equal to half the variogram range and the $Y(\mathbf{u}+\mathbf{h})$ values were extracted. Splitting the results into 50 quantiles and determining the mean and standard deviation of each quantile shows that the variance is homoscedastic for Gaussian data and the variance depends on the mean with lognormal data. For comparison, the analytical lines were also plotted on Figure 3.

Program

The slogsim program was adapted from the GSLIB sequential Gaussian Simulation algorithm, sgsim [6]. Three options are available for the type of simulation to be carried out and they are as follows:

Option 1	Transform a set of lognormal samples to Gaussian space and perform kriging and MCS, then back-transform to lognormal space. This is the
	standard/common approach. The limitation is that multiscale data are not easily handled.

- *Option 2* Perform direct kriging with the lognormal values with an adjusted variogram and do MCS without correcting the kriging variance. This is the published approach to DSS. The limitation is that heteroscedasticity/the proportional effect is not accounted for.
- *Option 3* Perform direct kriging on the lognormal values with an adjusted variogram and correct the kriging variance prior to MCS. This is the new approach that we are advocating in this paper. Multiscale data can be used in direct kriging and the proportional effect is explicitly accounted for.

Parameters are similar to sgsim; however, there are a few new ones to note:

START	OF PARAMETERS.	
Line	samples.dat	-file with data
1	1 2 0 4 0 0	-columns for X,Y,Z,vr,wt,sec.var.
2	-998.0 1.0e21	-trimming limits
3	1 0	-transform (0=no, 1=ves), dss (0=an,1=nai,2=peff)
4	slogsim.trn	- file for output trans table
5	0	- consider ref. dist (0=no, 1=yes)
6	histsmth.out	- file with ref. dist distribution
7	1 2	- columns
8	0.0 0.0	-Z mean and variance (if 0, determine from data)
9	0.0 0.0	-slope & intercept for prop effect fit, 0=not used
10	0.0 1500.0	-zmin,zmax(tail extrapolation)
11	1	-debugging level: 0,1,2,3
12	slogsim.dbg	-file for debugging output
13	slogsim.out	-file for simulation output
14	50	-number of realizations to generate
15	260 0.5 1.0	-nx,xmn,xsiz
16	300 0.5 1.0	-ny,ymn,ysiz
17	1 0.5 1.0	-nz,zmn,zsiz
18	69069	-random number seed
19	2 20	-min and max original data for sim
20	15	-number of simulated nodes to use
21	1	-assign data to nodes (O=no, 1=yes)
22	1 3	-multiple grid search (0=no, 1=yes), num
23	0	-maximum data per octant (0=not used)
24	150 75 1.0	-maximum search radii (hmax,hmin,vert)
25	165 0.0 0.0	-angles for search ellipsoid
26	50 50 1	-size of covariance lookup table
27	0 0.60 1.0	-ktype: 0=SK,1=OK,2=LVM,3=EXDR,4=COLC
28	/data/vdata.dat	-file with LVM, EXDR, or COLC variable

29 30 31 32 33	4 0 1 1	0.15 0.85	165 35	0.0 37	0.0	<pre>-column for secondary variable -Variogram option, 0=NScore Variogram, 1=Z-variogram -nst, nugget effect -it,cc,ang1,ang2,ang3 -a_hmax,a_hmin,a_vert</pre>
----------------------------	------------------	--------------	-----------	-----------	-----	--

The transform and dss options (Line 3) specify which form of simulation is to take place. If transform is set to 1, the input variable will be normal scored and sequential Gaussian simulation will be performed (dss should be set to zero in this case). To perform naïve direct lognormal simulation, transform should be set to zero and dss to one. To account for the proportional effect, dss must be set to 2 and transform to zero. On Line 8, the mean and variance of the data can be forced or determined from the input data. At this point, only a linear option for approximating the proportional effect is available (Line 9). If the proportional effect slope and intercept are input and the dss option from Line 3 is set to 2, kriging variance values will be corrected according to that linear equation. If the dss option is set to 2 and no slope/intercept parameters are entered, kriging is performed twice, once to acquire the Gaussian kriging variance and again for the estimate in original units. Line 30 is the last new parameter to describe the type of variogram that is input. If the variogram model was calculated in original space, this parameter must be set to zero; however, if the variogram model was calculated in original space, this parameter must be set to 1.

Output from the program is a gridded file containing simulated values similar to sgsim.

Examples

Tests were run on four data sets, one of which was simulated to be perfectly lognormal for comparison purposes. The data sets:

Data Set	Samples	Mean	Standard Deviation	Alpha	Beta
Analytical Lognormal	625	98.63	97.86	4.249	0.828
SIC Rainfall	467	184.24	112.26	5.058	0.562
Walker Lake, U	725	278.46	500.89	4.907	1.201
Walker Lake, V	725	277.64	249.23	5.331	0.769

The alpha and beta values were calculated from the actual sample values. These numbers are attained regardless of the type of distribution being input into the slogsim program and all conditional distributions resulting from kriging are assumed to be lognormal. Because of this, some initial testing was done on the sample sets to determine their lognormality, which will affect program output. As discussed, a random variable, $Z \mid z(\mathbf{u}) > 0$, is lognormal if the natural logarithm of $Z(\mathbf{u})$ is normally distributed. Since more robust statistical tests exist for normal distributions, the normality of ln(Z) was tested. The assumption made here is that data sets exhibiting more lognormal behavior will provide better results regarding mean, variance and variogram reproduction.

Three tests were carried out:

1. Fit error between the sample cumulative distribution function (cdf) and the best fitting analytical lognormal cdf. The fit error is the area between the two cdf's calculated using the Trapezoidal Rule for integration and is normalized by the width of the sample set.

- 2. Shapiro-Francia W` Test of Normality for large samples (5 to 5000 points), which is essentially the correlation between the standard normal cdf and the sample cdf of ln(Z).
- 3. A combined Skewness and Kurtosis measure defining how the sample probability distribution function (pdf) of ln(Z) deviates from normal [7]. For the standard normal distribution the skewness is zero, see Equation 11. The kurtosis magnitude is 3 so the equation is standardized by subtracting this value, see Equation 12.

$$Skewness = \frac{\sum_{i=1}^{N} (Y_i - \overline{Y})^3}{(N-1) \cdot \sigma^3}$$
(11)

$$Kurtosis = \frac{\sum_{i=1}^{N} (Y_i - \overline{Y})^4}{(N-1) \cdot \sigma^4} - 3$$
(12)

Where N is the number of samples, \overline{Y} is the mean and σ is the standard deviation. The error measure is the combined deviation these values indicated from standard normal where the deviations are defined by skewness normalized by the standard error of skewness (*ses*) and the kurtosis normalized by the standard error of kurtosis (*sek*), see Equations 13 and 14. Both have been normalized by the number of samples as well. The combined error measure is shown by Equation 15.

$$\sigma_{Skew} = \frac{|Skewness|}{\sqrt{6}} \qquad ses = \sqrt{\frac{6}{N}} \tag{13}$$

$$\sigma_{Kurt} = \frac{|Kurtosis|}{\sqrt{24}} \qquad sek = \sqrt{\frac{24}{N}} \tag{14}$$

$$\sigma_{Total} = \sqrt{\sigma_{Skew}^2 + \sigma_{Kurt}^2}$$
(15)

Lognormality test results for the analytical sample set (the best case) and the rainfall data can be found in Figure 6. The results are tabulated below:

Sample Set	Lognormal Fit Error	Shapiro-Francia W` Test	Skewness/Kurtosis Measure
Analytical Lognormal	0.004	0.997	0.071
SIC Rainfall	0.026	0.763	3.199
Walker Lake, U	0.017	0.868	0.603
Walker Lake, V	0.033	0.728	1.082

Variograms were modeled using the normal scores of the sample sets. The effects of lognormality on the variograms were checked by converting normal space variograms to lognormal space via Equation 8 and plotting them against the experimental variograms of the

samples in original units, see Figure 7. Original space experimental variograms of more lognormal data tend to match the analytical lognormal variogram model better.

Mean, variance and variogram reproduction were checked for 50 realizations of both the naïve and proportional-effect-corrected (Peff) forms of simulation. Maps of the mean and variance of the realization sets show homoscedastic variance for the naïve runs and reproduction of heteroscedasticity with runs accounting for the proportional effect, see Figure 4. Other Maps are located in Figure 8 (E-type means) and 9 (variances). Histograms of the mean and variance reproduction are shown in Figures 10 and 11 and summarized below. As expected, the reproduction of these two statistics was best with the analytical lognormal samples; however, even though the Walker Lake U variable seems more lognormal than V and the rainfall data, its mean and variance reproduction was worse.

	Original		Sim	ulated	Percent Error	
Data Set	Mean	Variance	Mean	Variance	Mean	Variance
Walker Lake, U	278.457	250542.3	234.97	211785.4	15.62	15.47
Walker Lake, V	277.638	62029.89	289.21	64664.23	4.17	4.25
SIC Rainfall	184.244	12576.05	158.5	15078.68	13.97	19.90
Analvtical Log	98.625	9560.733	100.25	9671.55	1.65	1.16

Mean and variance reproduction for naïve simulation

	Original		Sim	ulated	Percent Error	
Data Set	Mean	Variance	Mean	Variance	Mean	Variance
Walker Lake, U	278.457	250542.3	234.01	208120.4	15.96	16.93
Walker Lake, V	277.638	62029.89	288.69	67745.54	3.98	9.21
SIC Rainfall	184.244	12576.05	158.09	14179.26	14.20	12.75
Analytical Log	98.625	9560.733	100.19	9518.38	1.59	0.44

Mean and variance reproduction for Peff simulation

To check for variogram reproduction, a modified version of the GSLIB gam program was implemented over the sets of realizations. An average variogram was also calculated and compared visually with the input models. Variograms for both Naïve and Peff runs and for the major direction of anisotropy are shown, see Figure 12. Based on visual inspection, the input variogram models for both cases of direct lognormal simulation are reproduced well. The average variogram indicated by the dashed line tends to match the input model better with the proportional effect case. This is especially noticeable with the Walker Lake variables.

Conclusion

The direct simulation algorithm was written for use with lognormal data because it is analytically related to Gaussian data by the distribution and the variogram. It was also useful because the proportional effect is prominent and analytically defined. Many natural data sets exhibit the proportional effect, which direct kriging alone cannot reproduce. Introducing a variance correction into the kriging algorithm permits direct simulation where the proportional effect inherent in lognormal data is reproduced. Mean, variance, and variogram reproduction was acceptable for both the naïve and proportional-effect-corrected forms of direct simulation. Even though the naïve form indicated good variance reproduction, the variance was shown to be homoscedastic.

Having a direct simulation algorithm capable of handling the proportional effect will allow a move into data of various support volumes and population of unstructured grids. Advancing this

algorithm to one that can handle any input distribution and their inherent heteroscedastic features is sought after.

References

- 1 W.Xu and A.G. Journel, *DSSIM: A General Sequential Simulation Algorithm*, Stanford Centre for Reservoir Forecasting, Stanford University.
- A. Soares, *Direct sequential simulation and cosimulation*, Mathematical Geology 33 (8), 911–926, 2000
- J. Caers, *Adding local accuracy to direct sequential simulation*, Mathematical Geology 32 (7), 815–850, 2000a
- 4 B.Oz and C.V. Deutsch, *A Short Note on the Proportional Effect and Direct Sequential Simulation*, Centre for Computational Geostatistics, Report 4: 2001/2002.
- 5 A.G. Journel and Ch.J. Huijbregts, *Mining Geostatistics*, Academic Press, 1978, pp 570-573.
- 6 C.V. Deutsch and A.G. Journel, *Geostatistical Software Library and User's Guide*, Oxford University Press, second edition, 1998, pp 63-66.
- 7 M.R. Spiegel, *Theory and Problems of Probability and Statistics*, Schaum's Outline Series, McGraw-Hill Book Company, 1980

Appendix: Conditional Distributions for Direct Lognormal Kriging

To check if the local distributions at each location being estimated are lognormal, a simple kriging example was set up in Excel with 4 known data and 3 locations to be sequentially estimated:

120 -	120 Data Locations								
110 -	•	-1.5 20.3	-0.5 • 46.6						
		<mark>o</mark> Y*1							
≻ 100 -			o Y*3						
		<mark>o</mark> Y*2							
90 -	•	2.0		• 1.0					
		374		163					
80 -		100		110	 120				
			110	120					
Loc	cation	m_N		σ	.2 N				
Ŋ	/*1	-0.546	-0.546 0.4		62				
Ŋ	(*2	1.073		0.3	96				
Ŋ	(*3	0.267	1	0.4	44				

The data configuration for kriging is shown at the top and the kriging results prior to transformation to lognormal space are shown at the bottom. Solid bullets are known data and circles are the points to be estimated. The data values are also shown, both Gaussian (above) and lognormal (below).

To generate the local distributions corresponding to the global lognormal data with a mean and standard deviation of 100, a set of 199 quantiles was chosen ranging from 0.005 to 0.995 and the value corresponding to each for the local normal distributions was found. Using

$$Z(\mathbf{u}) = e^{\alpha + \beta \cdot Y(\mathbf{u})}$$

to transform the values to Z-space and plotting the results revealed that the local distributions are lognormal. To check if equations A1 and A2 (below) are correct, they were used to find the local alpha and beta values and then the LOGINV function in MS Excel was used to determine the corresponding value for each quantile. Both methods gave equal results, see below.



Starting with the equation for transforming Gaussian data, $Y(\mathbf{u}) \sim N(0, I)$, to normal X-space data, $X(\mathbf{u}) \rightarrow N(\alpha, \beta)$:

$$X(u) = \alpha_G + \beta_G \cdot Y(u)$$

Replacing Y(u) with the kriged mean and X(u) with the local mean of ln(Z), we get Equation A1, which relates the local mean in X-space to that in Gaussian space.

$$\alpha_L = \alpha_G + \beta_G \cdot m_N \tag{A1}$$

Where α_L is the local mean in X-space, α_G is the global mean of X(u), β_G is the global variance of X(u), and m_N is the kriged mean in Gaussian space. To derive Equation A2, the Equation for transforming Y-space values to X-space along with the equation defining the variance of a data set was used. The local normal kriging variance can be defined by the following equation:

$$\sigma_N^2 = \frac{1}{n} \sum_{i=1}^n (u_i - m_n)^2$$

Where σ_N^2 is the variance in Gaussian space, u_i is the value at location *i*, and m_N is the mean of all u_i , i=1...n. To determine the local variance of ln(Z) we need to know the values of u_{iX} and m_X that correspond to u_i and m_N in Gaussian space. Equation A1 can be used to perform this transformation:

$$u_{iX} = \alpha_G + \beta_G \cdot u_i$$
$$m_X = \alpha_G + \beta_G \cdot m_N$$

Substituting these into the equation for the variance in normal space (X-space), the local variance of ln(Z) can be solved for:

$$\beta_L^2 = \frac{1}{n} \sum_{i=1}^n (u_{iX} - m_X)^2$$

$$\beta_{L}^{2} = \frac{1}{n} \sum_{i=1}^{n} ((\alpha_{G} + \beta_{G} u_{i}) - (\alpha_{G} + \beta_{G} m_{N})^{2}$$

$$\beta_{L}^{2} = \frac{1}{n} \sum_{i=1}^{n} [\beta_{G} (u_{i} - m_{N})]^{2}$$

$$\beta_{L}^{2} = \frac{\beta_{G}^{2}}{n} \sum_{i=1}^{n} (u_{i} - m_{N})^{2}$$

$$\beta_{L}^{2} = \beta_{G}^{2} \cdot \sigma_{N}^{2}$$
(A2)

Where β_L^2 is the local variance of ln(Z), β_G^2 is the global variance of ln(Z), and σ_N^2 is the local normal variance in Y-space.



Figure 1: Normal and corresponding lognormal distributions



Figure 2: Variogram model used to generate the unconditional model. The Gaussian model is spherical with no nugget effect and a range of 32. The corresponding variogram of the lognormal variable is shown with the difference between the two functions.



Figure 3: Scatterplots of Gaussian data showing the variance is homoscedastic (left) and lognormal data displaying the proportional effect (right).



Figure 4: Variance of realizations for naïve and Peff simulation of the Walker Lake sample set.



Analytical Samples Lognormality

Figure 5: Lognormality of Analytical Lognormal Data.



SIC Rainfall Lognormality

Figure 6: Lognormality of Rainfall Data.



Figure 7: Normal scores and original units variograms.



Figure 8: E-type Mean Maps of the 50 Realizations.



Figure 9: Variance Maps of the 50 Realizations.



Figure 10: Mean and Variance Reproduction using Naïve Simulation Note: The darkest bar is that which is closest to the actual mean of the data.



Figure 11: Mean and Variance Reproduction using P-effect Simulation Note: The darkest bar is that which is closest to the actual mean of the data.



Figure 12: Variogram reproduction.