

On the classification and treatment of outliers in a spatial context: A Bayesian Updating approach

<http://dx.doi.org/10.1590/0370-44672021740003>

Victor Miguel Silva¹

<https://orcid.org/0000-0001-7502-9543>

¹AngloGold Ashanti - Avaliação de Recursos,
Mina Cuiabá, Sabará - Minas Gerais - Brasil.
victormiguelsilva@hotmail.com

Abstract

Checking and treating extreme values is commonplace in modelling workflows. The main methods to manage outliers may be categorized into graphical, Kriging- and simulation-based approaches. While graphical methods usually classify outliers from a global perspective, geostatistical methods evaluate outliers in a local context. Ordinary-Kriging based approaches are affected by conditional bias associated with the distribution tail(s), impacting on the correct classification of extreme values; the simulation method is based on the fact that geostatistical simulation is robust for outlier values. However, this approach ignores the interaction among outliers in the same neighborhood. The proposed approach considers that there are two values available at every sampled position, the sampled value and the conditional probability estimated from nearby data through cross-validation; the sampled value. Each value outside the user-defined threshold is classified as an outlier and is edited by merging the sampled and kriged value through Bayesian Updating. The proposed method is performed in normal-score units using Simple Kriging to (i) correctly estimate conditional distributions in the cross-validation step; (ii) avoid conditional bias; and (iii) minimize the outlier influence on experimental-variogram modelling. The proposed method is compared to three other widely used methods in a case study of a gold deposit. The proposed method substantially improved the local accuracy and reduced the number of misclassified blocks of a reference model.

Keywords: outliers; extreme values; Bayesian updating; geostatistics.

1. Introduction

The presence of outliers dramatically impacts many steps of data analysis, such as exploratory data analysis, experimental variogram development and distribution fitting. Next, in the modelling step, just a few outliers may overestimate the variability and the estimated grade of a large region, over or under-estimating the global grade and/or the volume of ore and waste at a given cut-off. Barnett and Lewis (1984, pp. 32) give the following definition of an outlier: “An outlier in a set of data is an observation that appears to be inconsistent with the remainder of that set of data”. In general, we may differentiate the methodologies for outlier classification based on: how they define the datum context, how they measure the consistence among data; and how a datum designated as an outlier is managed.

In the field of geostatistics, the context is given by the complete dataset or subsets of statistically similar groups, such as samples with the same

rock type, sampling support and type, located in the local neighbourhood, among other criteria. Graphical methods define the context under the assumption that the inlier-data distribution is represented by a given curve or theoretical distribution (Sichel, 1966; Srivastava, 2001). Geostatistical-based methods may be used prior to estimates to edit outliers based on their relationship with neighborhood data (Hawkins and Cressie 1984; Costa, 2003; Maleki et al. 2014) or the estimates may be directly performed using robust algorithms (Arik and Kim 1992; Machado et al. 2012; Babakhani 2014). Further details of these methods are provided in the following section.

The proposed approach considers that there are available two sets of reliable data at every sampled position: (i) the conditional probability estimated from nearby data through cross-validation; and (ii) the sample value associated to a inferred minimal variance that assures a degree of

overlapping area, controlled by the C value between it and the kriged conditional distribution. Both distributions are merged through Bayesian Update, an inference method derived from the theorem of Bayes (Bayes, 1763). The cross-validation is performed using normal-score data and Simple Kriging (SK) to correctly estimate the conditional distributions and to avoid the impact of conditional bias in the outlier classification.

This article is structured as follows: We review the advantages and limitations of methods presented in literature; and next, we present a new approach to edit outliers by integrating the datum value with the kriged distribution at their position, weighted by their associated uncertainty. Then, a gold deposit composed of four systems of mineralized quartz-carbonate is used to illustrate the application of the proposed approach and its comparison with other methodologies. A conclusion follows.

1.1 Managing outliers

Outliers are commonly considered as the result of local anomalies that are not representative of the entire domain or region to be modelled, or resulting from human error, improper sample preparation and/or assay problems. The first step

in outlier management is to fix or discard data associated with errors. A second step is to rethink how we have grouped data and if there are clues of the existence of unrecognized subpopulations (Srivastava, 2001).

Commonly it is not possible to deter-

mine the cause of all inconsistencies and outliers, and therefore, the values designated as outliers are expected to combine non-detected inconsistencies and geological anomalies. The following methods show alternatives to classify and handle outliers.

1.2 Graphical solutions

The graphical approaches are the simplest and frequently used in the mining industry to define outliers. The premise is that the true distribution of the data is represented by a given curve (David, 1988) and the inlier data are expected to follow a fairly constant

behaviour along this curve. Probability plots are widely used to visually define extreme values that detach from common, random variations associated with the data distribution or its curve (Figure 1a). The visual inspection is supported in most sophisticated approaches by an

analytical comparison of the data with a chosen distribution, such as two or three-parameter lognormally distribution (Sichel, 1966; Figure 1b). In all cases, extreme values distant from the expected behaviour are classified as suspicious values.

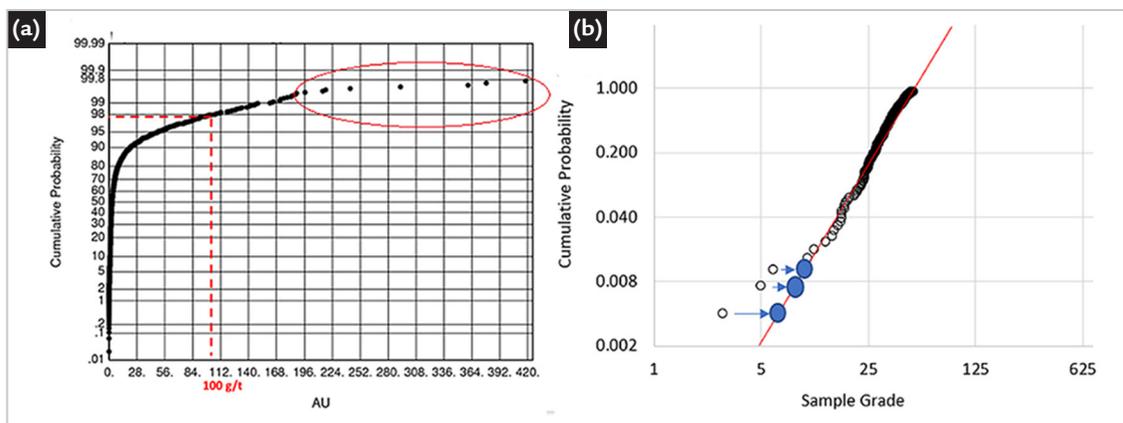


Figure 1 – Cumulative probability plot in a) original units, where the high-grade outliers are within the red ellipse. b) lognormal units, where low-grade outliers (blue points) are adjusted to the fitted distribution (red line).

Figure 2a presents the cutting curve (Roscoe, 1996), representing the sensitivity of the average grade to the different capping values. The capping value is defined as the grade near the

inflection point before the plateau where the curve stabilizes. Figure 2b shows the data's cumulative coefficient of variation (CV) in descending order. The value near the point in which there is a pronounced

increase in the CV is defined as the capping value. The CV curve is proposed as one step of the workflow developed by Parker (1991), which is further discussed in the "Kriging-based solutions" section.

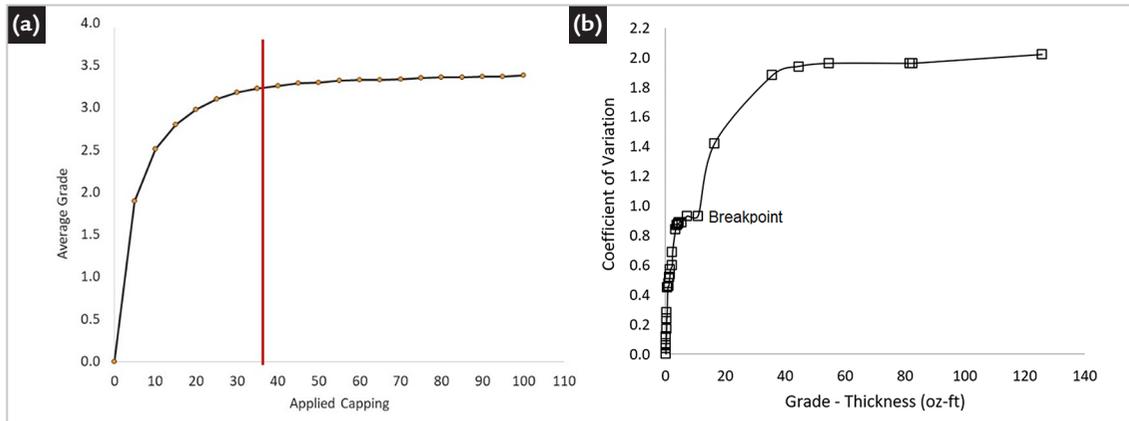


Figure 2 - a) Cutting curve between the applied capping and the average grade, where the vertical-red line indicates where the curve stabilizes; b) Cumulative coefficient of variation of the data in descending order against the "grade x thickness" product. Adapted from Parker, 1991.

Analyzing outliers in a global context is appropriate to manage sparsely sampled models if the available data is not enough to estimate their local distribution.

When there are enough data to estimate reliable local distributions, it is relevant to consider that spatial data value may be or not be an outlier in different positions

along the deposit. We use an exhaustive unidimensional dataset to illustrate the problem of classifying local outliers from a global perspective (Figure 3).

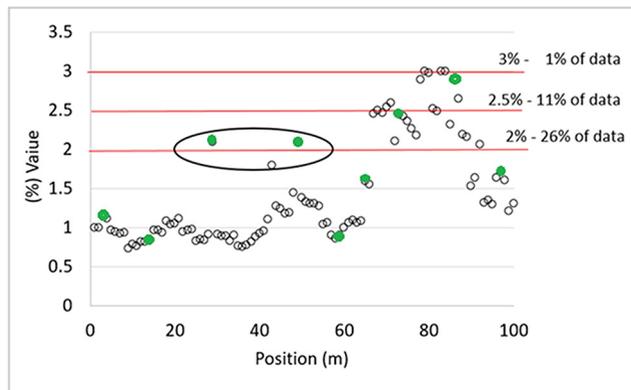


Figure 3 - Scheme between three different capping values (2%, 2.5% and 3%) and their relationship with three local outliers (within the black ellipse) in an exhaustive dataset. Green points mimic sparsely sampled data extracted from the exhaustive dataset.

In Figure 3 the green circles represent data sampled from an exhaustive dataset. The exhaustive dataset may be divided into one domain with grade between 0.7-1.5% in the interval [0, 0.65 m] and a second domain [65 m, 100 m] where the

grades range between 1.2-3.1%. Three outliers may be observed in the low-grade domain within the black ellipse. In this data setting, three applied capping values removed inlier values and left part (or all) outliers unedited. While domaining

is a clear solution when we analyze the exhaustive dataset, this solution is not clear if only the sparsely sampled grid is available (green circles). The following section presents geostatistical methods to check outliers in their spatial context.

1.3 Kriging-based solutions

Robust estimators are the first Kriging-based approaches proposed to manage extreme values, such as Lognormal Kriging (LK) (Journel, 1980) and multiple Indicator Kriging (mIK) (Journel, 1982). More recently the field parametric geostatistics was proposed as a solution to manage outliers (Machado *et al.*, 2012). The LK is extremely sensitive to the sill of the variogram and can overestimate the

values when the lognormality assumption does not hold. The mIK commonly requires to be post-processed to fix order relation violations, and support correction, both can be an additional source of bias in the final estimate (Costa, 2003). Parker (1991) used the CV curve to classify outliers (Figure 2b), being these values adjusted to a lognormal distribution and their spatial continuity individually esti-

mated by Indicator Kriging (IK; Journel, 1982). Maleki *et al.* (2014) proposed to jointly cokriging the truncated-grade data and the weighted indicator of grade above the defined threshold. Fourie *et al.* (2019) proposed to restrict the kriging weights attributed to outliers proportionally to their probability of occurrence in the data distribution.

The Robust kriging (RoK) is pro-

posed to minimize the influence of outliers using its neighboring data (Cressie and Hawkins, 1980; Hawkins and Cressie, 1984). Costa (2003) proposed an approach based on RoK, where each datum

value $z(\mathbf{u}_i)$ is compared to a robust median-weighted estimate $z^{wm}(\mathbf{u}_i)$ at its position \mathbf{u}_i . The estimator weights are computed from the cross-validation step performed by Ordinary Kriging (Matheron, 1963). The

method may check low and/or high-value outliers. If the deviation between $z(\mathbf{u}_i)$ and the robust estimate $z^{wm}(\mathbf{u}_i)$ is larger than the user-defined threshold (Equation 1), $z(\mathbf{u}_i)$ is replaced by an edited value $z^e(\mathbf{u}_i)$

$$z^e(\mathbf{u}_i) \begin{cases} z^{wm}(\mathbf{u}_i) + c\sigma_{ok}(\mathbf{u}_i), & \text{if } z(\mathbf{u}_i) - z^{wm}(\mathbf{u}_i) > c\sigma_{ok}(\mathbf{u}_i) \\ z(\mathbf{u}_i), & \text{if } |z(\mathbf{u}_i) - z^{wm}(\mathbf{u}_i)| \leq c\sigma_{ok}(\mathbf{u}_i) \\ z^{wm}(\mathbf{u}_i) - c\sigma_{ok}(\mathbf{u}_i), & \text{if } z(\mathbf{u}_i) - z^{wm}(\mathbf{u}_i) < -c\sigma_{ok}(\mathbf{u}_i) \end{cases} \quad (1)$$

Where $\sigma_{ok}(\mathbf{u}_i)$ is the kriging standard deviation and the C value is a user-defined constant, which controls the amount of editing of each outlier. Further details about the value C are presented in the Proposed Method section (Step iii).

The RoK based approach considers the datum in its local context, it does not require any additional variogram

1.4 Simulation-based solution

Babakhani (2014) proposed geostatistical simulation to calibrate a cutting level for estimation based on the fact that geostatistical simulation is a robust technique with respect to outlier values, where extreme values are not exaggerated by the local weighting of estimation, since the simulated realizations are constrained to reproduce the data distribution and spatial continuity. This approach was further developed by Chiquini and Deutsch (2017), who proposed the following steps to use simulated values to calibrate the cutting level to be used in estimates:

(i) The values to be checked are usu-

ally chosen from the distribution tail(s), ignoring local outliers that do not stand out as extreme values in global terms;

(ii) The volume of influence of each datum value to be estimated is delimited. The volume includes all blocks that the kriging weight attributed to the datum of interest is above a determined value;

(iii) The mean and uncertainty of each volume of influence is established by geostatistical simulation; and

(iv) Each datum value is adjusted until the estimated quantity of metal is similar to the simulated values in its volume of influence.

ally incorrect and may not converge to the underlying real variogram (Babakhani 2014). Secondly, the conditional bias associated with Ordinary Kriging estimates affects the distribution tail(s), which is the area of interest in the outlier classification. These problems are addressed by the simulation-based solution discussed in the next section.

The simulation-based method solves two issues related to kriging-based solutions: its conditional-bias is null, and the use of an experimental variogram fitted on normal score units mitigates the influence of extreme values, clustering and the proportional effect (Deutsch and Kumara, 2017). The method, however, ignores the local interaction among overlapping outliers, and therefore, it is appropriate only to check a small number of samples whose volumes of influence do not overlap each other. Next, we present a novel methodology considering the limitations and advances of the methods discussed above.

2. Proposed method

The proposed approach relies on the understanding that there are two data sets at each sampled position: (i) the datum value with its associated uncertainty $N(y(\mathbf{u}_i), \sigma_{error}^2(\mathbf{u}_i))$; and (ii) the conditional probability $N(y^*(\mathbf{u}_i), \sigma_{sk}^2(\mathbf{u}_i))$ estimated by Simple Kriging (SK) using cross-validation from neighbor data. The most probable inlier value at \mathbf{u}_i results from merging both distributions through Bayesian Updating.

The values are transformed to normal-score units and those estimated by SK are used in the step of cross-validation to avoid conditional bias because they allow that all marginal and conditional distributions be fully defined by their mean and variance. The proposed method, hereafter referred to as BU, consists of the following steps:

(i) **Transforming values:** Original data values $z(\mathbf{u}_i)$ are transformed into normal-scored values to $y(\mathbf{u}_i)$ to meet part of the

assumptions required by multi-gaussian normal equations (equivalent to Simple Kriging). Moreover, adjusting the experimental variogram to transformed values partially mitigates the influence of extreme values, clustering and the proportional effect (Deutsch and Kumara, 2017);

(ii) **Cross-validation:** Cross-validation is performed at all sampled positions using SK. The cross-validation consists of temporarily removing one-by-one each sample or the entire drillhole and using the remaining samples to estimate $N(y^*(\mathbf{u}_i), \sigma_{sk}^2(\mathbf{u}_i))$ at the node of the removed sample. The process is repeated at each sampled position \mathbf{u}_i . The decision of removing each sample or the entire drillhole in each step depends on the type of mineralization and whether the most probable origin of the outliers is individually associated with each sample or simultaneously associated with all

samples along the drillhole.

The use of Simple Kriging is necessary because its regression slope is exactly 1. In the case of ordinary kriging where the Lagrange Multiplier is used, it is typically less than one, indicating a conditional bias which can increase the number of values misclassified as outliers;

(iii) **Defining the data threshold:** The BU method relies on the assumption that the larger is the deviation between the data value $y(\mathbf{u}_i)$ and the estimate $N(y^*(\mathbf{u}_i), \sigma_{sk}^2(\mathbf{u}_i))$, the larger is the chance of the sampled value being an outlier. All values outside the threshold interval $[-c\sigma_{sk}(\mathbf{u}_i), c\sigma_{sk}(\mathbf{u}_i)]$ are assumed as outliers, where the C value is a user-defined parameter (Cressie and Hawkins, 1980; Costa, 2003). Figure 4 shows these parameters (note that the location index \mathbf{u}_i was dropped from the notation).

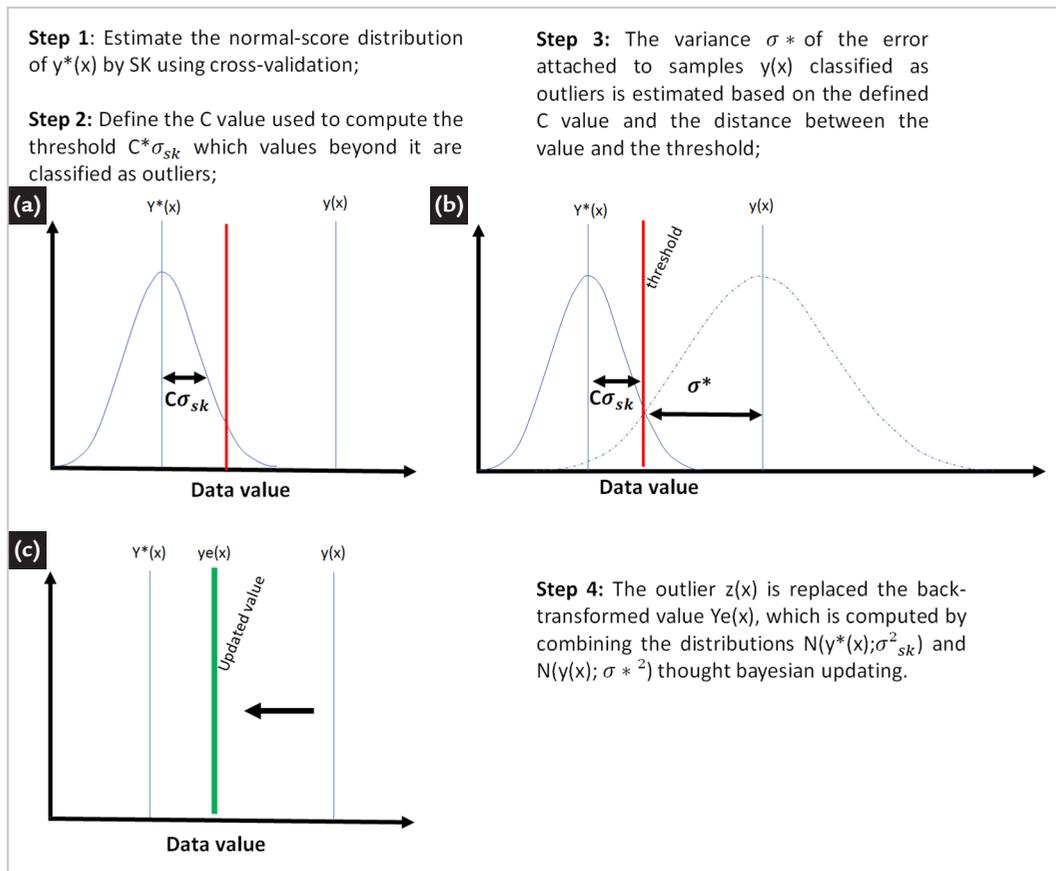


Figure 4 – Scheme of the proposed workflow. a) The conditional distribution kriged by Simple Kriging and the user-defined parameter C are used to define the threshold whereby all values beyond are classified as outliers; b) The variance associated with the datum value at the kriged position is estimated (right-side distribution); c) The sampled value classified as outlier and kriged values are merged by Bayesian Updating. The dataset with treated values is used to estimate the model of interest.

The bigger the C value, the smaller the change made to the original datum value. Cressie and Hawkins (1980) cite that the C value in the 1.5–2.5 interval is recommended and that values of C below 1.0 proved to excessively smooth out the outlier effect. However, Hawkins highlights that there are no fixed criteria

to define outliers and these values should be analyzed and adjusted for each case. Costa (2003) suggests from a practitioner perspective that the calibrated C value should not modify the declustered global mean by more than 5%.

(iv) *Estimating the uncertainty of each outlier:* The uncertainty $\sigma_{error}^2(\mathbf{u}_i)$ as-

$$\sigma_{error}(\mathbf{u}_i) = \frac{|y(\mathbf{u}_i) - C\sigma_{sk}(\mathbf{u}_i)|}{C} \tag{2}$$

Equation 2 computes the minimal variance $\sigma_{error}^2(\mathbf{u}_i)$ associated to the outlier under analysis $y(\mathbf{u}_i)$ that assures a degree of overlapping area between $N(y(\mathbf{u}_i), \sigma_{error}^2)$

(\mathbf{u}_i) , and the kriged distribution $N(y^*(\mathbf{u}_i), \sigma_{sk}^2(\mathbf{u}_i))$, given the user-defined C value and its threshold $[-C\sigma_{sk}(\mathbf{u}_i), C\sigma_{sk}(\mathbf{u}_i)]$. The larger this overlapping area, the larger the

sociated with each outlier $y(\mathbf{u}_i)$ is assumed as resulting from the combination of microscale variation, human error, improper sample preparation, and/or assay and the minimal error intrinsically associated with any sampling process. In the case of an outlier in the upper tail (high-grade outliers), the error $\sigma_{error}(\mathbf{u}_i)$ is given by

probability that the data under analysis is an inlier value. Figure 5 shows the relationship between three different distributions with the threshold $C\sigma_{sk}(\mathbf{u}_i)$.

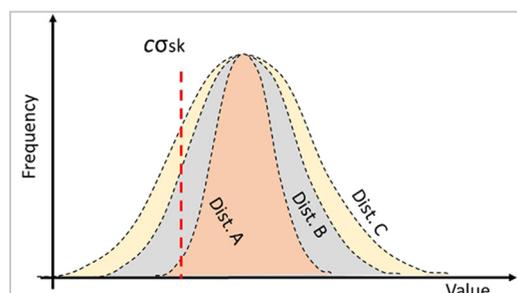


Figure 5 – Value threshold (vertical-dashed line) based on the user-defined C parameter and the distribution estimated by Simple Kriging through cross-validation. The higher the distribution area of the lower tail of the distributions A, B and C below the threshold limit, the higher the probability that the extreme value is inlier.

The practitioner may choose to treat only anomalous values that are above or below the kriged value at its neighborhood, or both. In these cases, the only difference is the signal associated with the C value, being possible to define simultaneously both thresholds. In general, it is

recommended to set C as positive when managing outliers of a positively skewed distribution and use -C if the distribution is negatively skewed. Next, the value classified as an outlier and the kriged distribution at its position are combined into an updated value;

$$y_{bu}^e(\mathbf{u}_i) = \frac{y^*(\mathbf{u}_i) \sigma_{error}^2(\mathbf{u}_i) + y(\mathbf{u}_i) \sigma_{sk}^2(\mathbf{u}_i)}{\sigma_{error}^2(\mathbf{u}_i) - \sigma_{sk}^2(\mathbf{u}_i) \sigma_{error}^2(\mathbf{u}_i) + \sigma_{sk}^2(\mathbf{u}_i)} \quad (3)$$

Where equation 3 is the Bayesian Update equation derived from the Theorem of Bayes whereby both available data combinations are weighted by their associated uncertainty.

3. Case study

3.1 Geological settings

The study area comprises a typical quartz alteration mineralization with subordinated sulphides assemblage which is located in the Archean Greenstone Belt in the *Quadrilátero Ferrífero* region, an important Brazilian metallogenic province. The orebody of interest is composed of mineralized quartz veins

(vi) *The updated values are back transformed to original units.*

Next, the proposed approach is applied to a gold deposit whose gold distribution is highly skewed. For the

(v) *Estimating the updated value y_{bu}^e* : Each value $y(\mathbf{u}_i)$ beyond the defined threshold are classified as outlier value and replaced by the new value $y_{bu}^e(\mathbf{u}_i)$. Bayesian Updating equation is used to combine $N(y(\mathbf{u}_i), \sigma_{error}^2(\mathbf{u}_i))$ and $N(y^*(\mathbf{u}_i), \sigma_{sk}^2(\mathbf{u}_i))$ in order to generate the updated value $y_{bu}^e(\mathbf{u}_i)$:

sake of confidentiality, the case study used a synthetic dataset mimicking the geological, statistical and geostatistical characteristics of the real gold deposit under study.

hosted in the stratigraphic footwall of the banded iron ores that compose the main orebodies of this deposit. The orebody of interest is the quartz “vein-associated” composed of four quartz-carbonate vein systems (Figure 6): V1 mineralized system is a fault-fill type associated with the principal gold-bearing fluid input stage.

The mineralized V2 veins are controlled by the foliation and subdivided according to internal crystal orientation. The extensional array V3 veins are also controlled by the foliation, whereas the extensional, breccia-style V4 system has veins associated with flanking structures (Vitorino, 2017).

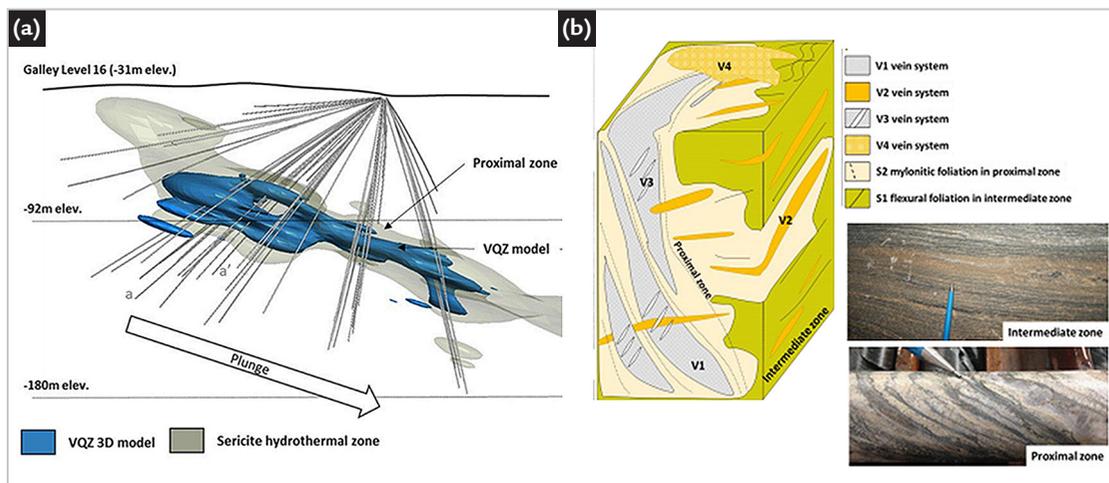


Figure 6 - Geological model from the mineralized zone. a) looking parallel to the plunge. b) schematic organizations showing the gold grade variability and organizations of the quartz-carbonate veins systems. Adapted from Vitorino (2017).

The orebody complexity and lack of additional information make domaining the model into individual veins systems difficult and very subjective.

3.2 Case-study methodology

The comparison among methods is based on a reference model assumed as ground truth, which was estimated using a synthetic short-term/grade-control dataset composed of 342 samples with average length of 70 cm (Figure 7b; 7c).

All the proximal zone, therefore, is modelled together and a grade shell above 1 g/t. is delimited. The grade magnitude order of each vein system is

The methods to be compared used a subset extracted from the complete dataset, composed of 176 samples with an average spacing of 30 m was. Samples in both the test and ground-truth datasets were given equal weight in estimates.

different, making the chosen orebody a good case study to evaluate the capacity of different methods to classify local outliers.

The criteria to compare methods and to choose their parameters must consider the model purpose. In this case, we prioritize the model capacity of correctly classifying blocks as ore or waste from estimates using the long-term sampling grid.

3.3 Dataset and method comparison

The BU method is compared with Robust Kriging (RoK; Costa, 2003); Capping calibrated by simulation (SIM; Babakhani, 2014; Chiquini and

Deutsch, 2017) and Ordinary Kriging performed using the raw data without outlier treatment. These methods used the long-term sampling grid composed

of 176 samples which were extracted from the complete short-term/grade-control dataset composed of 342 samples (Figure 7b; 7c).

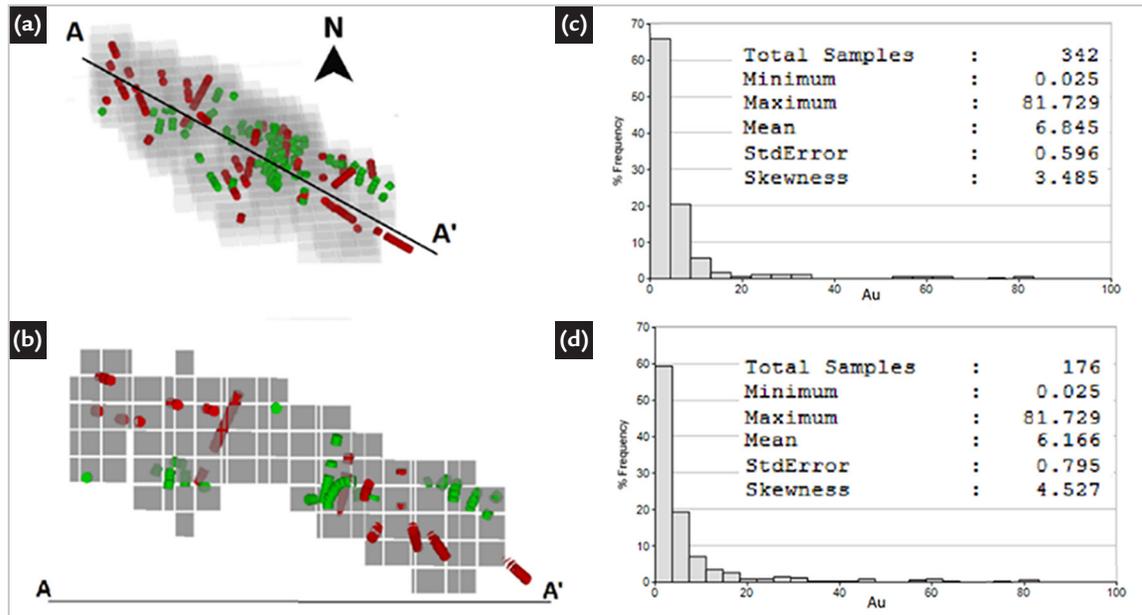


Figure 7 - 3D geological model from orebody mineralized zone. a) the plan of view is parallel to the plunge; b) Cross-section along AA'. Histograms of declustered data with the c) complete dataset with production data and d) test dataset.

The variogram $\gamma_{orig}(\mathbf{h})$ and $\gamma_{ns}(\mathbf{h})$ were modelled for the complete dataset in original and normal-score units. Both models indicated directional

anisotropy, being the major continuity along the mineralization plunge and the minor along the vertical direction. The models are composed of the nug-

get effect and one spherical structural with directional ranges indicated in parentheses (Equation 4-5).

$$\gamma_{orig}(\mathbf{h}) = 7 + 20 \cdot \text{Sph}\left(\frac{Plun.}{40m} \quad \frac{Strike}{30m} \quad \frac{Vert.}{8m}\right) \quad (4)$$

$$\gamma_{ns}(\mathbf{h}) = 0.21 + 0.79 \cdot \text{Sph}\left(\frac{Plun.}{45m} \quad \frac{Strike}{34m} \quad \frac{Vert.}{9m}\right) \quad (5)$$

The search ellipse distances of all methodologies were equal to the directional ranges of their variogram models, with a minimum of 8 and a maximum of

24 samples. The C value was set 2 for RoK and BU, which in both cases is the largest rounded-up value that does not modify the global declustered mean more than

approximately 5%.

Figure 8 shows the relationship between the C value and the accuracy of estimates using the BU approach.

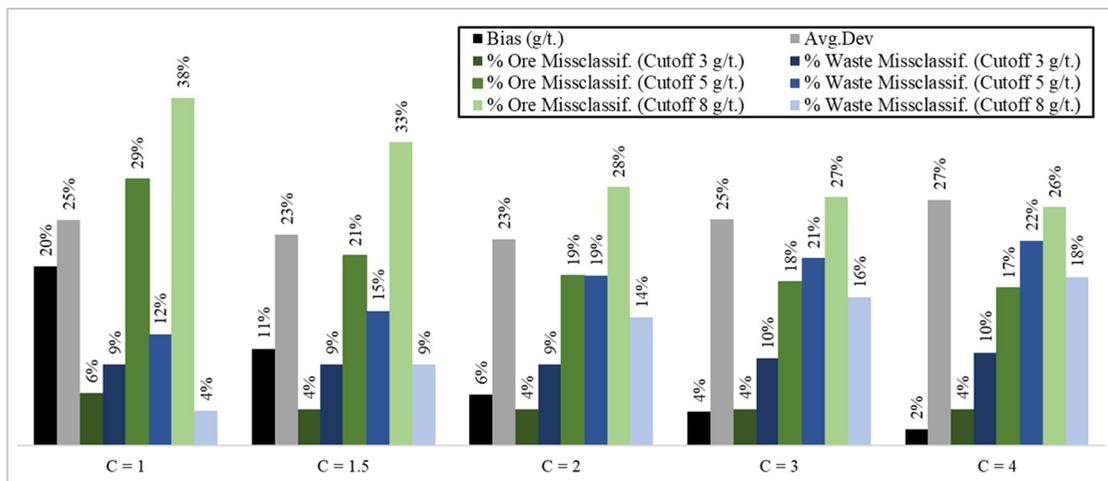


Figure 8 - Bias, rate of waste and ore misclassification and the average grade deviation between the reference block model and the model estimated using BU approach for different values of C.

We may observe in the range [1, 4] that the largest is the C value, the lowest is the misclassification and the average deviation between the estimated and reference grade of each block. However, increasing C also increases the bias compared to the declustered average grade.

Figure 9a shows the scatterplot between the samples in normal-score units and the kriged values performed by cross-validation. The decision of removing each sample or the entire drillhole in the step

of Cross-validation of the BU approach needs to consider the probable origin of the anomalies to be treated in each specific problem. In this case study, each sample was removed one-by-one. The author considers that anomalous high-grade values to be treated in this deposit are associated with centimetric intervals that are smaller than the sample length of 70 cm and thus, probably influence a single sample and not the entire drillhole. Moreover, human error, improper sample

preparation, and/or assay are also related with a single sample.

Figure 9b shows the difference between the originally sampled values and the new values edited by the proposed approach with the C value set to 2. The correction was performed in normal-score units, using the original data transformed to a gaussian distribution with declustered mean 0 and unit variance, and then back-transformed to original units after the estimation.

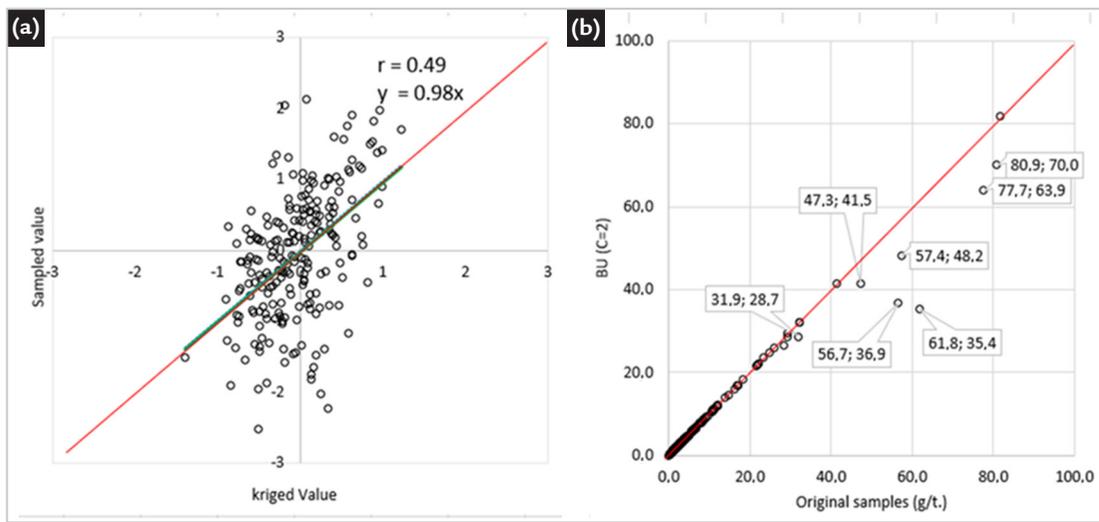


Figure 9 – Scatterplot for a) cross-validation performed using SK to estimate the normal-score values, and b) the effect of editing original-unit values using the BU approach and C set to 2.

Considering the model of reference, Table 1 shows the average deviation, correlation and the rate of ore and waste misclassification of Ordinary Kriging per-

formed using the raw data, BU, and RoK for three cut-off grades. It is important to highlight that the average grade of the reference dataset is higher than the declus-

tered grade of the test dataset, respectively 6.84 g/t. and 6.16 g/t. Therefore, it is expected that the evaluated models present a negative bias against the reference model.

Table 1 – Comparison of estimates of Robust Kriging (RoK), Bayesian-Updating approach (BU),

Capping calibration by simulation (SIM), and Ordinary Kriging using the raw dataset without any treatment (OK).

All using the same test dataset composed of 176 samples with declustered mean of 6.16 g/t. The methods are compared to a ground-truth model estimated using the complete dataset with 342 samples with declustered mean of 6.84 g/t.

METHOD	REFERENCE	ROK	BU	SIM	OK
Avg. Grade (g/t.)	6.64	5.77	5.74	5.82	5.83
Avg. Deviation	-	41%	31%	34%	45%
Correl. Coefficient r	1	0.41	0.54	0.50	0.39
Cutoff – 3 g/t.					
% Ore misclassified	-	14%	4%	5%	9%
% Waste misclassified	-	6%	9%	10%	11%
% Misclassified	-	11%	6%	7%	15%
Cutoff – 5 g/t.					
% Ore misclassified	-	34%	21%	17%	18%
% Waste misclassified	-	24%	15%	21%	22%
% Misclassified	-	30%	18%	19%	21%
Cutoff – 8 g/t.					
% Ore misclassified	-	63%	60%	55%	55%
% Waste misclassified	-	12%	8%	9%	14%
% Misclassified	-	26%	22%	22%	24%

The BU approach is the more efficient method in limiting the influence of outliers from a local perspective. The better local accuracy is indicated by the smaller misclassification rate for all tested cut-offs (Table 1), the

higher coefficient of correlation, and the lower standard error of regression (Figure 10). The smaller spread results in bigger stability among models estimated with new samples added from the long to the grade-control

model. The cloud of points above the reference X=Y line indicates blocks influenced by inlier values incorrectly classified as outliers; the cloud below the line results from values misclassified as inliers.

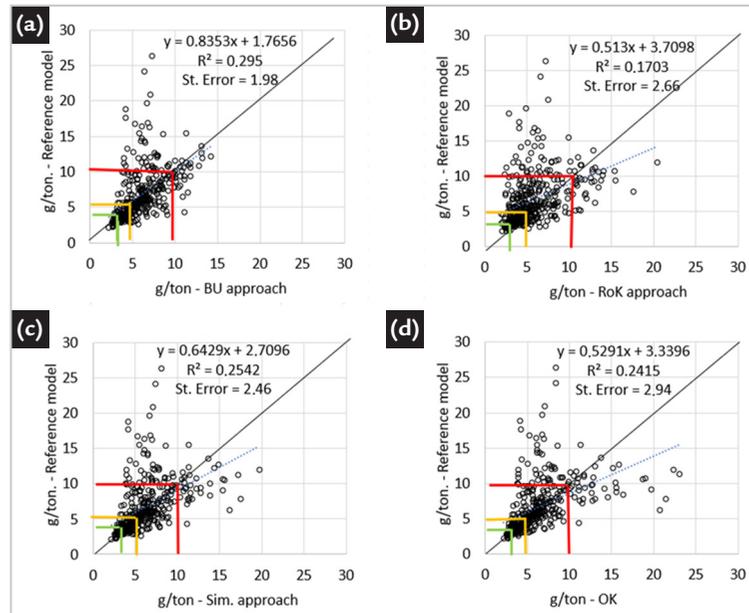


Figure 10 – Scatter plot between the reference model and the same blocks estimated by a) Bayesian Updating based approach; b) Robust Kriging; c) Capping calibration by simulation; and d) Ordinary kriging using the raw dataset without any treatment. The green, yellow and red lines indicate the cut-off grade of respectively 3, 5 and 8 g/t.

The grade-tonnage curves are shown for the overall assessment of recoverable resources from a global perspective at different cut-offs. The behavior of the curves

of the estimated and reference models are presented in Figure 11. While the behavior of all compared methodologies are very similar below 7.5 g/t., the proposed

approach underperforms the other methodologies above 7.5 g/t., being specially biased above 10 g/t., which corresponds to less than 5% of the total tonnage.

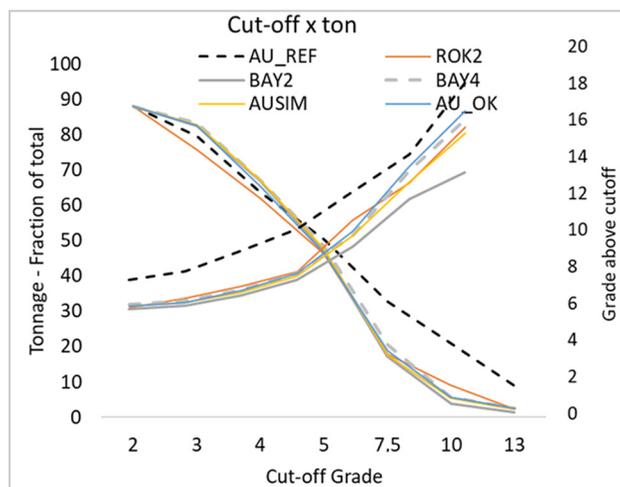


Figure 11 – Grade tonnage curves (descending curves) and mean grade per cut-off (ascending curves). The reference curves (black-dashed) were obtained from estimates using the complete dataset.

In this case study, we applied different methods to treat outliers for a highly skewed deposit of gold. The BU approach estimates have the best the local, followed by the simulation and RoK approaches. BU estimates above 7.5 g/t, however, are globally biased. These results highlight that

the C value must be set considering the model application.

Higher C value leads to biased estimates that are a concern from a global perspective but improves the local accuracy and reduces the misclassification rate, which is of paramount importance to decide if the block should be sent to

the waste pile or the processing plant. If the main purpose were the global statistics, BU approach with C set to 4 would be a better parameter because it approximates the estimates to the reference grade x ton curve (Figure 11) but, as expected, would worsen local accuracy (Figure 8).

4. Conclusions

The proposed approach addresses problems associated to other geostatistical solutions, such as the experimental-variogram modelling using the data in normal-score units, reducing the influence of extreme values; the conditional bias is null as Simple Kriging has a slope of regression of exactly one; and the method evaluates all available data as a function

of their similarity with nearby values.

The capacity of the proposed approach for correcting, classifying and treating outliers was tested in a highly skewed gold deposit composed of many generations of mineralization with abrupt changes from high to low-grade areas. The proposed approach outperformed the other compared methods

available in literature in the case which the local accuracy is prioritized, improving the misclassification rate and the deviation between the reference and the estimated blocks. We recommend analysing in future studies the performance of the proposed method with other data distributions and variable levels of skewness.

References

- ARIK, A.; KIM, Y. C. Outlier restricted kriging: a new kriging algorithm for handling of outlier high grade data in ore reserve estimation. *In: INTERNATIONAL SYMPOSIUM ON THE APPLICATION OF COMPUTERS AND OPERATIONS RESEARCH IN THE MINERAL INDUSTRY*, 23., 1992, Littleton, Colorado. *Proceedings [...]*. New York: Society of Mining Engineers of AIME, 1992. p. 181-187.
- BABAKHANI, M. *Geostatistical modeling in presence of extreme values*. 2014. 85 f. Thesis (Master of Science in Mining Engineering) - University of Alberta, Edmonton, Canada, 2014.
- BARNETT, V.; LEWIS, T. *Outliers in statistical data*. 2. nd. Chichester: John Wiley & Sons, 1984. 386p.
- BAYES, P. R. T. An essay towards solving a problem in the doctrine of chances. By the late Rev. Mr. Bayes, F. R. S. communicated by Mr. Price, in a letter to John Canton, A. M. F. R. S. *Philosophical Transactions of the Royal Society of London*, v. 53, p. 370-418, 1763.
- CHIQUINI, A. P.; DEUTSCH, C. V. A simulation approach to determine the cutting level. *In: Deutsch, J. L. (ed.). Geostatistics lessons. [S. l.: s. n.]*, 2017. Available at: <http://www.geostatisticslessons.com/lessons/simulationcutting>. Accessed: 26 Jan. 2021.
- COSTA, J. Reducing the impact of outliers in ore reserves estimation. *Mathematical Geology*, v. 35, n.3, p. 323-345, 2003.
- CRESSIE, N.; HAWKINS, D. M. Robust estimation of the variogram. *Mathematical Geology*, v. 12, n. 2, p. 115-125, 1980.
- DAVID, M. *Handbook of applied advanced geostatistical ore reserve estimation*. Amsterdam: Elsevier Science, 1988. 216p.
- DEUTSCH, C. V.; KUMARA, P. Transforming a variogram of normal scores to original units. *In: Deutsch, J. L. (ed.). Geostatistics lessons. [S. l.: s. n.]*, 2017. Available at: <http://www.geostatisticslessons.com/lessons/convertnsvariograms>. Accessed: 26 Jan. 2021.
- FOURIE, A.; MORGAN, C.; MINNITT, R. C. A. Limiting the influence of extreme grades in ordinary kriged estimates. *Journal of the Southern African Institute of Mining and Metallurgy*, v. 119, n. 4, p. 391-401, 2019.
- HAWKINS, D. M.; CRESSIE, N. Robust kriging: a proposal. *Journal of the International Association for Mathematical Geology*, v. 16, n. 1, p. 3-18, 1984.
- JOURNAL, A. G. The lognormal approach to predicting local distributions of selective mining unit grades. *Mathematical Geology*, v. 12, n. 4, p. 285-303, 1980.
- JOURNAL, A. G. The indicator approach to estimation of spatial distributions. *In: INTERNATIONAL SYMPOSIUM ON THE APPLICATION OF COMPUTERS AND OPERATIONS RESEARCH IN THE MINERAL INDUSTRY*, 17., 1982, Golden, Colorado. *Proceedings [...]*. New York: AIME, 1982. p. 793-806.
- MACHADO, R. S.; ARMONY, M.; COSTA, J. F. C. L. Field Parametric Geostatistics: a rigorous theory to solve problems of highly skewed distributions. *In: ABRAHAMSEN, P.; HAUGE, R.; KOLBJØRNSEN, O. (ed.). Geostatistics Oslo 2012*. Dordrecht: Springer, 2012. p. 383-395.
- MALEKI, M.; MADANI, N.; EMERY, X. Capping and kriging grades with long-tailed distributions. *Journal of the Southern African Institute of Mining and Metallurgy*, v. 114, n. 3, p. 255-263, 2014.
- MATHERON, G. Principles of geostatistics. *Economic Geology*, v.58, n. 8, p. 1246-1266, 1963.
- PARKER, H. Statistical treatment of outlier data in epithermal gold deposit reserve estimation: *Mathematical Geology*, v. 23, n. 2, p. 175-199, 1991.
- ROSCOE, W. E. Cutting curves for grade estimation and grade control in gold mines. *In: ANNUAL GENERAL MEETING, CANADIAN INSTITUTE OF MINING, METALLURGY AND PETROLEUM*, 98., 1996, Edmonton, Alberta. *Proceedings [...]*. [S. l.: s. n.], 1996. 8p.
- SICHEL, H. The estimation of means and associated confidence limits for small samples from Lognormal populations. *Journal of the Southern African Institute of Mining and Metallurgy*, special symposium volume, p. 106-123, 1966. (Symposium on Mathematical Statistics and Computer Applications in Ore Valuation).
- SRIVASTAVA, M. *Outliers: a guide for data analysts and interpreters on how to evaluate unexpected high values*. Vancouver, British Columbia: [s. n.], 2001. 4 p. (Contaminated Sites Statistical Applications Guidance Document

N. 12–8).

VITORINO, L. A. *Mineralização aurífera associada aos veios quartzo-carbonáticos hospedados na unidade máfica basal da jazida Cuiabá, Greenstone belt Rio das Velhas, Minas Gerais, Brasil*. 2017. 98 f. Dissertação (Mestrado em Geologia) – Instituto de Geociências, Universidade Federal de Minas Gerais, Belo Horizonte, 2017.

Received: 23 December 2020 - Accepted: 3 May 2021.



All content of the journal, except where identified, is licensed under a Creative Commons attribution-type BY.