Measuring sampling fairness and geochemical consensus for blasthole assays within grade-block mining units

A PREPRINT

Raymond Leung

Rio Tinto Centre for Mine Automation Australian Centre for Field Robotics (ACFR) The University of Sydney Sydney, NSW 2006 raymond.leung@sydney.edu.au

June 21, 2022

ABSTRACT

In the mining industry, current grade control practices lack a standardised framework that can assess the reliability of average grade estimates computed for selective mining units (otherwise known as grade-blocks) within a mining bench. This article describes two measures that can quantify sampling fairness and geochemical consensus. Concretely, sampling fairness considers spatial factors such as sampling density and bias in the spatial distribution of blastholes whereas geochemical consensus considers the agreement between the assay samples within a grade-block. Geochemical disparity is measured using a robust distance estimator and a masking formula that takes into account the proportion of outliers and magnitude of differences observed above a threshold. The efficacy of the consensus measure is demonstrated through validation experiments. The results confirm the MCD robust estimator can breakdown when the fraction of outliers exceeds (n-k-1)/(2n). For $k \ge 2$ variables and a sample size $n \ge 10$, this typically leads to an underestimation of the true extent and impact of outliers when they exceed 40%. An extension based on split-sequence analysis is proposed to overcome this limitation. The method is evaluated using production data from an open-pit iron ore deposit. An open-source implementation of the proposed algorithms will be available on github.

Keywords Selective mining units \cdot grade-blocks \cdot sampling fairness \cdot geochemical consensus \cdot robust estimator \cdot high-breakdown

1 Introduction

In open-pit iron ore mining, a bench is partitioned into selective mining units (SMU) to facilitate grade control and targeted excavation [1]. The SMU concept originates in geostatistical estimation and is often understood as the smallest volume of material on which ore and waste classification is determined. In the context of banded iron-ore formation (BIF) hosted iron ore deposits [2] in the Hamersley Province of Western Australia, it is a common practice to divide a bench into SMUs known as *grade-blocks*. These polygons vary in shape and size from 100 m² to >10,000 m². The intent is to decompose a region in such a way that each grade-block maps to one destination. This generally ensures the material excavated from a grade-block will be transported to a stockpile or waste dump as appropriate. Ideally, the composition of each grade-block is homogeneous, so the average grade computed from blasthole assays collected from the grade-block would produce a representative value. In practice, this average grade may not be representative if (i) the blasthole sampling is sparse or highly irregular; (ii) significant variability or even incongruous measurements in the assay samples are ignored. In essence, the *reliability* of the average grade-block value depends on two things: fairness in the spatial sampling process, and the geochemical consensus among the assays within a grade-block. Within the mining industry, there is currently no standard framework for expressing these ideas. This article describes two objective measures that allow sampling fairness and geochemical consensus to be quantified.

It allows the appropriateness of a particular grade-block configuration in a pit to be evaluated to indicate whether further refinement (e.g. reshaping or splitting certain grade-blocks) would improve ore-waste delineation. It is also worth noting that grade prediction block models are often evaluated against grade-block values which serve as the ground-truth. An example of this is shown in [3]. Thus, incorporating a reliability measure for individual grade-blocks—for instance, discounting or putting less emphasis on grade-blocks where data is lacking or incongruous—would enable one to draw strength from more confident data and provide a more informed comparison during R_2 model reconciliation [4]. Our motivation is to develop a capability where unrepresentative assay mean estimates can be highlighted in the affected grade-blocks within a mining bench.

The concepts of sampling fairness and geochemical consensus are further explained with pictures in Fig. 1. Fig. 1a shows an example where dense, uniform blasthole spacing yields high spatial confidence. In contrast, the distribution of the blastholes is sparse and irregular in Fig. 1b, this yields lower spatial confidence. In terms of geochemistry, Fig. 1c shows an example where consistent composition would produce a high consensus score. For simplicity, only one chemical component is shown. Fig. 1d depicts a situation where the assay samples within the block are incongruous. In this case, the proposed method would produce a low consensus score. With the foundation now firmly established, the methodology will be described in the next section.



2 Formulation: measures for sampling fairness and geochemical consensus

Two measures are formulated to characterise the spatial distribution and geochemical consensus among blasthole assay samples within a grade-block. The approach targets the following observations which help explain the main factors that erode confidence in the grade estimates.

- Low sampling density few blastholes are drilled in some grade-blocks.
- Non-uniform sample distribution blasthole spacing is highly irregular.
- *Inadequate grade-block coverage* sampling is biased or sparse, often there is no information (holes not drilled) for a large portion of the grade-block.

These **spatial factors** (limited knowledge or confidence in the reference values themselves) can sometimes explain poor predictive performance of a grade model. Apart from undersampling, **chemical factors** also influence the reliability of grade estimates:

- Fraction of outliers grade-block may contain unrepresentative samples that deviate from the mean. For
 instance, most of the samples might be classified as mineralised (high grade), except a few might be classified
 as unmineralised (low grade). This may result from inaccurate localisation of the mineralisation boundary.
- *Magnitude of difference* the extent to which the assays disagree within the samples also matters. Of special interest are any exceptional fluctuations within the data beyond what is considered normal.

2.1 Spatial factors

The spatial factors considered are concerned with fairness in the grade-block blasthole sampling pattern.

2.1.1 Sampling entropy

Fig. 2 shows the blasthole distribution of several grade-blocks with varying degree of bias and sparsity. To measure how uniformly distributed the samples are within each grade-block \mathcal{G}_i , a kD-tree is built using the N_i blasthole coordinates, $\{\mathbf{h}_{i,j}\}_{1 \le j \le N_i}$. Separately, K_i i.i.d. particles are synthesized from a uniform distribution to cover the entire grade-block. The scope of each sample may be described by Voronoi cells. This is equivalent to finding the closest blasthole sample associated with each of the random particles in a nearest neighbour search. Using $p_{i,j}$ to denote the fraction of particles associated with blasthole j in grade-block i, the entropy is given by

$$f_{\text{entropy}}(i) = -\frac{\sum_{j=1}^{N_i} p_{i,j} \log_2 p_{i,j}}{\max\{\log_2 N_i, 1\}} \in [0, 1]$$
(1)

The entropy approaches 1 when all $p_j \equiv p_{i,j}$ are equally probable. This happens when samples are evenly spread, roughly speaking, when holes are drilled on a regular (rectangular or hexagonal) grid. The entropy, $f_{entropy}(i)$, is reduced progressively as blasthole samples become more non-uniformly spaced. The sampling entropy may be interpreted as a measure of fairness or regularity in the sampling pattern.



Figure 2: Spatial factors relating to the distribution of blasthole samples within grade-blocks. Top: blastholes shown as red dots, i.i.d. particles rendered in distinct colours in each Voronoi cell. Bottom: sensitivity of the $f_{coverage}$ metric to grade-blocks with varying degree of coverage.

2.1.2 Grade-block coverage

Querying the kD-tree also returns $\{d_k^{\min}\}_{1 \le k \le K_i}$, the distance between the k^{th} particle, \mathbf{p}_k , and its closest hole. The inverse cumulative distribution functions (icdf) of d_k^{\min} are shown in Fig. 2(bottom). The range of influence is defined

as the average blasthole spacing observed in a grade-block. This is represented by $R_i^{(influence)} = \text{median}\lambda_j$ where $\lambda_j = \min_k \{\|\mathbf{h}_{i,j} - \mathbf{h}_{i,k}\|\}_{j \neq k}$ is the minimum separation for sample *j* from another blasthole. The proposed measure for grade-block coverage is then given by

$$f_{\text{coverage}} = 1 - \frac{\left| \left\{ \mathbf{p}_k \mid d_k^{\min} > R_i^{(\text{influence})} \right\} \right|}{K_i} \tag{2}$$

where the fraction in expression (2) represents the proportion of points outside the range of influence of the assayed blastholes [see blue shaded area in Fig. 2].

2.1.3 Sample density

The inverse of density is *sparsity*. Sparsity may be defined as the number of holes drilled within a grade-block divided by its area. In the context of open-pit mining at Pilbara iron ore deposits, it is more convenient to measure density in units of $100m^2$ instead of $1m^2$. With $\lambda_{\text{density}} \stackrel{\text{def}}{=} 100 \times (\text{number of samples within})/(\text{grade-block area } [m^2])$, the proposed measure for sampling density is given by

$$f_{\text{sample-density}} = 1 - \beta \cdot \alpha^{\lambda_{\text{density}}} \text{ with } \beta = 0.3, \alpha = 0.1$$
 (3)

where the penalty term $\alpha^{\lambda_{\text{density}}}$ is restricted to [0, 1] (the smaller the sampling density exponent, the higher the penalty) and the baseline β is set at 0.3 (i.e. the minimum value of $f_{\text{sample-density}}$ is ≥ 0.7) to ensure the strength of all spatial factors (f_{entropy} , f_{coverage} and $f_{\text{sample-density}}$) are comparable.

2.1.4 Spatial confidence

The proposed spatial confidence measure is computed as

$$f_{\text{spatial-confidence}} = (f_{\text{coverage}})^{1/2} \cdot f_{\text{entropy}} \cdot f_{\text{sample-density}}$$
(4)

The spatial confidence map for a single bench in a Pilbara iron ore mine is shown in Fig. 3. The grade-blocks with the lowest value generally correspond to waste or low-grade blocks since there is little incentive to sample an area well if the material is destined for a waste dump.



Figure 3: Map of spatial confidence for grade-blocks in a single bench, $f_{\text{spatial-confidence}}$

2.2 Chemical factors

Factors that contribute to grade-block reliability are often multi-factorial. Suboptimal sampling does not necessarily mean the assay measurements will be problematic. Conversely, fair sampling also does not guarantee no conflict between the assay samples. Therefore, the assessment is incomplete without examining variability within the grade-block assay samples. Robust statistical estimators provide a useful starting point in identifying a subset of trustworthy samples within a grade-block. Measuring deviation using Mahalanobis distance, MVE (minimum volume ellipsoid), MCD (minimum covariance determinant) and pruning MST (minimum spanning tree) are some of the candidates for finding supportive samples that minimise data dispersion.

The MCD technique [5] possesses several desirable attributes, such as affine equivariance¹, high break-down point and the existence of fast algorithms. It is therefore used in the current formulation. MCD provides a robust estimation of the location and scatter of the data by minimising the determinant of the covariance matrix. This procedure excludes outliers that unduly influence the raw estimates. Its goal is equivalent to minimising the differential entropy for Gaussian distributed data. However, this interpretation is not valid for compositional data such as assay measurements. In bivariate analysis, when the coordinates corresponding to a pair of chemical components are plotted, the points are not Gaussian or even symmetrically distributed [6]. Instead, the data lies in a simplex governed by Aitchison geometry [6]. In order to work with Euclidean distances, MCD is applied to isometric log-ratio (ILR) transformed data $\{ilr(\mathbf{z}_{i,j})\}_{1\leq j\leq N_i}$, where $\mathbf{z}_{i,j}$ denotes the %wt concentration of chemical components in an assay sample j within grade-block \mathcal{G}_i .

Geochemical assays often include a full array of elements. In the case of iron ore, these elements may comprise [Fe, SiO₂, Al₂O₃, P, LOI, Mn, MgO, S, CaO, TiO₂]. Usually, only a subset of these are considered important for analysis. Given a particular subcomposition, say $\mathbf{z} = [Fe, SiO_2, Al_2O_3]$, $\mathbf{c}_{i,j} = C[\mathbf{z}] \in \mathbb{R}^m$ (here, m = 3) simply represents the subcomposition vector subject to closure. If we drop the subscripts *i* and *j* which index the *j*th assay sample in grade-block \mathcal{G}_i ,

$$\mathbf{c} = \mathcal{C}[\mathbf{z}] = \frac{\mathbf{z}}{\sum_{k=1}^{m} z_k}$$
(5)

The ILR is subsequently computed using the centred log-ratio (CLR) transformation and the Helmert sub-matrix described by Tsagris et al. [7]. Concretely, the CLR transformation is defined as

$$\operatorname{clr}(\mathbf{c}) = \mathbf{y} = [y_1, ..., y_m] \in \mathbb{R}^m$$
, where $y_i = \log \frac{c_i}{\sqrt[m]{\prod_{k=1}^m c_k}} = \log c_i - \frac{1}{m} \sum_{k=1}^m \log c_k$ (6)

The ILR transformation may be calculated as

$$\operatorname{ilr}(\mathbf{c}) = \mathbf{H}\mathbf{y} \tag{7}$$

where $\mathbf{H} \in \mathbb{R}^{(m-1) \times m}$ corresponds to a row-normalised version of M—the Helmert matrix [8]

$$\mathbf{M} = \begin{pmatrix} 1 & -1 & 0 & 0 & \dots & 0 \\ 1 & 1 & -2 & 0 & \dots & 0 \\ \vdots & & \ddots & & \vdots \\ 1 & 1 & \dots & 1 & -(m-2) & 0 \\ 1 & 1 & \dots & 1 & 1 & -(m-1) \end{pmatrix}.$$
 (8)

Row k in the orthonormal matrix **H** is given by

$$\mathbf{H}[k,:] = \left[\underbrace{\frac{1}{\sqrt{k^2 + k}}, \dots, \frac{1}{\sqrt{k^2 + k}}}_{k \text{ terms}}, \frac{-k}{\sqrt{k^2 + k}}\right] \text{ for } 1 \le k < m.$$
(9)

These operations achieve scale invariance and subcompositional coherence [6] which address the issue of spurious correlation when analysis is performed on raw data. It also means interpretations based on this and other subcompositions should yield similar results that are consistent with results obtained using the whole composition.

¹It is scale and rotation invariant. The latter means the analysis is insensitive to the orientation of the correlation structure.

2.2.1 Fraction of outliers

The first factor considered is the proportion of unrepresentative samples present within a grade-block. This is estimated via MCD which computes the *robust distance* of each ilr-transformed assay sample from the centre of the cluster. Writing $\mathbf{x}_{i,j} = \mathrm{ilr}(\mathbf{c}_{i,j}) \in \mathbb{R}^{m-1}$, the robust distance of sample j in grade-block i is given by $d_{i,j} = \sqrt{(\mathbf{x}_{i,j} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1}(\mathbf{x}_{i,j} - \boldsymbol{\mu}_i)}$.² The main point to note is that $d_{i,j}^2$ is χ^2 distributed. Hence, samples with $d_{i,j} > \sqrt{\chi_{p,\nu}^2}$ are identified as non-representative samples. For p = 0.975 and $\nu = 2$, $\sqrt{\chi_{p,\nu}^2} \approx 2.716203$. Accordingly, the fraction of outliers, $\eta_{\text{outliers}} \ge 0$, is estimated as

$$\eta_{\text{outliers}} = \frac{\left|\left\{j \mid d_{i,j} > \sqrt{\chi_{p,\nu}^2}\right\}\right|}{N_i} \tag{10}$$

2.2.2 Magnitude of difference

To evaluate the extent of disagreement among the assay samples in a grade-block, we first compute the *geometric mean of the robust distance for the outliers*:

$$d_{\text{outliers}}^{\text{gmean}} = \exp\left(\frac{1}{|S_{i,\text{outliers}}|} \sum_{j \in S_{i,\text{outliers}}} \ln d_{i,j}\right)$$
(11)

Next, $\sqrt{\chi^2_{p,\nu}}$ is used to measure how far the distortion is above the threshold. The masking function is given by

$$d_{\text{outliers}}^{\text{unmasked}} = \min\left[\max\left[\log_{10}\left(\frac{d_{\text{outliers}}^{\text{gmean}}}{\sqrt{\chi_{p,\nu}^2}}\right), 0\right], 1\right]$$
(12)

which limits the range to [0,1]. Values close to 0 (respectively, 1) are interpreted as trivial (respectively, significant) differences.

2.2.3 Geochemical consensus

The overall consensus among samples within a grade-block is given by

$$f_{\rm consensus} = \left(1 - \eta_{\rm outliers}\right)^{d_{\rm outliers}^{\rm unmasked}} \tag{13}$$

This power-law expression is simple yet elegant. It describes a decaying curve that decreases from 1 to $1 - \eta_{\text{outliers}}$ as the unmasked distortion, $d_{\text{outliers}}^{\text{unmasked}}$, increases. It represents a family of curves each parameterised by η_{outliers} . Within a grade-block, for a given proportion of outliers (say $\eta_{\text{outliers}} = 0.3$), the level of consensus is lower-bounded by $1 - \eta_{\text{outliers}} (= 0.7)$. This occurs in the worst-case scenario where the disagreement between the assay samples is extreme. On the other hand, if the samples barely disagree, i.e., only slightly above the threshold, $\operatorname{sqrt}(\chi^2_{p,\nu})$, then the consensus will be pushed back toward 1. This adaptive behaviour gives due regard to the level of conflict between the grade-block samples. Beyond η_{outliers} , the deciding factor is the unmasked distortion, $d_{\text{outliers}}^{\text{unmasked}}$.

2.2.4 Small sample size adjustment

The MCD method may not work effectively when the sample size (n) is small. According to Hubert and Debruyne [5], it requires at least $n \ge 2m$ samples, where m denotes the number of features (variables). To avoid the curse of dimensionality, a heuristic is used to estimate consensus when there are no more than six samples (using n > 3m as a rule of thumb). The median \mathbf{m}_i , median absolute deviation $\boldsymbol{\sigma}_i$, standard scores $\mathbf{z}_{i,j}$ and coefficient of variation \mathbf{r}_i are computed for assay measurements, $\mathbf{c}_{i,j} \in \mathbb{R}^m$, taken from grade-block \mathcal{G}_i as follows:

$$\boldsymbol{\sigma}_{i} = \text{median}\left(|\mathbf{c}_{i,j} - \mathbf{m}_{i}|\right), \quad \mathbf{z}_{i,j} = \frac{\mathbf{c}_{i,j} - \mathbf{m}_{i}}{\boldsymbol{\sigma}_{i}}, \quad \mathbf{r}_{i} = \frac{\boldsymbol{\sigma}_{i}}{\mathbf{m}_{i}} \in \mathbb{R}^{m}$$
 (14)

The fraction of outliers, $\hat{\eta}_{\text{outliers}}$, and conflict scores, $\hat{d}_{\text{conflict}}$, are computed as

$$\hat{\eta}_{\text{outliers}} = \frac{\left|\{j \mid \mathbf{w}^T \mid \mathbf{z}_{i,j} \mid > \lambda_i\}\right| + 1}{N_i + 1} \tag{15}$$

$$\hat{d}_{\text{conflict}} = \min\left\{\ln(1 + \lambda_i \mathbf{w}^T \mathbf{r}_i), 1\right\} \in [0, 1]$$
(16)

 $^{{}^{2}\}mu_{i} \in \mathbb{R}^{m-1}$ and $\Sigma_{i} \in \mathbb{R}^{m-1 \times m-1}$ represent robust mean and covariance estimates obtained from $h \leq J$ supporting samples.

where the weight vector $\mathbf{w}^T = [0.5, 0.325, 0.175]$ corresponds to [Fe, SiO₂, Al₂O₃] and the function $\lambda_i = 3 \times \left(\frac{2}{\pi} \tan^{-1} \sqrt{N_i}\right)^4$.³ Geochemical consensus is estimated as

$$\hat{f}_{\text{consensus}} = (1 - \hat{\eta}_{\text{outliers}})^{d_{\text{conflict}}} \text{ when } N_i \le 6$$
 (17)

2.2.5 Example

Table 1(left) shows the composition of assay samples within two grade-blocks: HGB6 and WH7 before ilrtransformation. Looking at the Fe column, the geochemical variation is relatively small for the HGB6 grade-block. In contrast, there is significant variation in Fe (likewise for SiO_2 and Al_2O_3) for the WH7 grade-block which contains a mixture of shale and low-grade mineralised samples. Using MCD, the robust distances computed for each sample are sorted and displayed on the right-hand side of Table 1. In particular, distances that exceed the $\chi^2_{p,\nu}$ critical value of 2.7162 are printed in bold. As intuition would suggest, there is a higher proportion of outliers in WH7 than HGB6. Furthermore, the magnitude of the outliers relative to $\sqrt{\chi^2_{p,\nu}}$ are also much higher. The quantities of interest from (10), (11), (12) and (13) are computed and shown in Table 2. The values agree with intuition: a low consensus score for WH7 reflects the significant disparity between the assay measurements. Fig. 4a shows a map of geochemical consensus for the same bench used in Fig. 3. Grade-block samples that lack consensus are coloured in lighter shades. To verify some of the results, Fig. 4b highlights three grade-blocks (VHA37, VHA38 and W44) with particularly low $\hat{f}_{consensus}$ values of 0.718, 0.756 and 0.733, respectively. Fig. 5 uses voronoi tessellation to show the geochemical variation of Fe, SiO₂ and Al₂O₃ within these grade-blocks. Each voronoi cell is centered at the blastholes where assay samples were taken. The colour intensity of each cell is proportional to the concentration of the respective chemical components. Significant variation can be seen in all three grade-blocks. It is clear that their compositions are not homogeneous; thus a mean estimate based on the sample average (in Fe, say) will not be representive of the entire grade-block. The $\hat{f}_{\text{consensus}}$ score allows these instances to be automatically detected. A grade-block with low geochemical consensus may be split if its samples can be clustered. In the case of VHA37 (see Fig. 5a-5c), the horizontal band running through the middle of the grade-block may be segmented to form a new grade-block to reduce the grade sample variance.

| Table 1: MCD | robust dist | ance shows | the level o | of disparity | among sam | ples in two | grade-blocks |
|--------------|-------------|------------|-------------|--------------|-----------|-------------|--------------|
| | | | | | 6 | | 0 |

| col | mposition | of assay sa | mples wit | hin grade-b | lock | sorted sample robu | ust distance |
|--------|-----------|-------------|-----------|------------------------------|-----------|--|----------------------|
| HC | GB6 (N = | 15) | Ī | $\overline{\text{WH7}}(N=2)$ | 25) | HGB6 | WH7 |
| Fe | SiO_2 | Al_2O_3 | Fe | SiO_2 | Al_2O_3 | $d_{\mathrm{HGB6},j}$ | $d_{\mathrm{WH7},j}$ |
| 61.516 | 3.2701 | 4.1705 | 43.044 | 14.4198 | 13.4280 | 0.4346 | 0.5609 |
| 61.667 | 4.5334 | 3.6154 | 47.794 | 12.61 | 10.1063 | 0.6883 | 0.5620 |
| 61.947 | 3.7471 | 3.4418 | 28.535 | 23.8188 | 20.7988 | 0.7163 | 0.5825 |
| 61.968 | 2.5144 | 2.2479 | 44.452 | 11.2035 | 12.1828 | 0.8486 | 0.6590 |
| 62.364 | 4.3374 | 3.6550 | 39.832 | 40.0663 | 1.8570 | 0.9689 | 0.6960 |
| 62.390 | 3.7766 | 3.0398 | 37.105 | 43.296 | 1.0949 | 1.0969 | 0.7854 |
| 62.433 | 4.2752 | 3.4948 | 48.456 | 12.3453 | 10.3742 | 1.3537 | 0.8309 |
| 63.046 | 3.2867 | 2.9829 | 48.019 | 11.6673 | 10.6452 | 1.5195 | 1.0866 |
| 64.003 | 3.1212 | 2.6463 | 54.677 | 6.2498 | 7.8491 | 1.7372 | 1.2544 |
| 64.879 | 1.9877 | 1.6556 | 41.319 | 15.9348 | 14.7593 | 2.1607 | 1.2547 |
| 64.966 | 2.8495 | 2.0632 | 62.100 | 2.9384 | 2.9736 | 2.4765 | 1.4712 |
| 56.955 | 8.2475 | 6.4583 | 61.059 | 3.3467 | 3.2409 | 4.0303 | 1.5934 |
| 57.052 | 4.4025 | 4.7908 | 39.752 | 40.0334 | 1.2384 | 4.0819 | 1.6942 |
| 58.023 | 3.6744 | 5.0130 | 61.538 | 7.9433 | 1.5573 | 4.2308 | 2.2061 |
| 59.882 | 4.1615 | 5.3099 | 36.024 | 42.8717 | 1.1932 | 4.7751 | 2.4218 |
| | | | 45.053 | 30.9200 | 1.7437 | | 2.4836 |
| | | | 56.202 | 14.6952 | 1.8480 | | 5.3103 |
| | | | 38.866 | 40.9796 | 1.1927 | | 10.8539 |
| | | | 36.003 | 44.8525 | 0.9443 | | 15.7432 |
| | | | 36.653 | 44.0994 | 0.9393 | | 15.8015 |
| | | | 39.105 | 40.4532 | 1.2416 | | 16.0474 |
| | | | 50.281 | 24.9943 | 1.5467 | | 18.0692 |
| | | | 58.421 | 10.2444 | 3.0714 | | 19.7316 |
| | | | 36.193 | 45.7901 | 1.0260 | | 21.2168 |
| | | | 36.228 | 43.8815 | 1.1215 | $\sqrt{\chi^2_{p,\nu}} \approx 2.7162$ | 29.1046 |

³For N_i in [2,3,4,5,6], it generates λ_i of 0.4104, 0.5926, 0.7404, 0.8626, 0.9658, respectively.

| | | HGB6 | WH7 |
|--------------------------------------|---|-------------------------|-----------------------|
| fraction of outliers | $\eta_{\rm outliers}$ | $\frac{4}{15} = 0.2666$ | $\frac{9}{25} = 0.36$ |
| robust distance geometric mean | $d_{\text{outliers}}^{\text{gmean}}$ | 4.2697 | 15.4670 |
| robust distance unmasked distortion | $d_{\text{outliers}}^{\text{unmasked}}$ | 0.196442 | 0.755444 |
| geochemistry consensus among samples | $f_{\rm consensus}$ | 0.9408 | 0.7138 |

Table 2: Geochemical analysis: statistics computed for the same two grade-blocks



Figure 4: Map of geochemical consensus for grade-blocks in a single bench, $\hat{f}_{\text{consensus}}$

2.3 Overall reliability

Finally, an overall measure for the reliability of the grade-block estimates may be computed from $\hat{f}_{\text{spatial-confidence}}$ and $\hat{f}_{\text{consensus}}$. This is simply given by

$$\hat{f}_{\text{reliability}} = \hat{f}_{\text{spatial-confidence}} \cdot \hat{f}_{\text{consensus}}$$
 (18)

A map of $\hat{f}_{\text{reliability}}$ for the same bench is shown in Fig. 6.



Figure 5: Assay sample disparity within grade-blocks with low geochemical consensus scores



Figure 6: Map of overall reliability for grade-blocks in a single bench, $f_{\text{reliability}}$

3 Validation

This section examines the efficacy of the proposed measure for geochemical consensus. The objective is to show $\hat{f}_{\text{consensus}}$ scales proportionally with the fraction and magnitude of outliers contained in the assay samples; that it generally behaves as a decreasing function of η_{outliers} and $d_{\text{outliers}}^{\text{unmasked}}$. Since these parameters are unknown, test data is synthesized so that these attributes can be tightly controlled. This process is informed by trends observed in the blasthole assay data endemic to a Pilbara iron ore mine. Fig. 7a shows the data points ($\mathbf{c} \in \mathbb{R}^3$) in a ternary diagram. Using vector quantisation, point density is more clearly depicted in Fig. 7b using a heat map. The contour plot in Fig. 7c shows the data is concentrated along two axes angled at $\theta_0 = 304.7^\circ$ and $\theta_1 = 273.9^\circ$ from the apex.



Figure 7: Trends observed in the compositional data endemic to a Pilbara iron ore mine site



Figure 8: Synthesis of validation data: parameters governing the sampling procedure

Parameters that govern data synthesis are shown in Fig. 8. After selecting the anchor points $(\mathbf{p}_0^{\text{start}}, \mathbf{p}_0^{\text{end}}, \mathbf{p}_1^{\text{start}})$ and $\mathbf{p}_1^{\text{end}}$, the principal components that correspond to the dotted lines are given by φ_0 and φ_1 where $\varphi_m = \mathbf{p}_m^{\text{end}} - \mathbf{p}_m^{\text{start}}$. The main variables are η , δ and ν , where $\eta \in [0, 1]$ denotes the fraction of points drawn from the lower ellipse, $\delta \in [0, 1]$ adjusts the distance between two given ellipses, and the noise parameter ν controls the semi-major length of the sampling ellipse. In the experiments, η varies from 0 to 1 in increments of 0.0625, δ varies from 0.05 to 0.95 in increments of 0.05, and ν is fixed at 0.025. For each configuration (m, δ, η) , θ_m controls the orientation of the sampling ellipses, $(1 - \eta)N$ random points are selected inside the upper ellipse at $\mathbf{p}_m^{\text{start}}$, the remaining ηN points are selected from the lower ellipse at $\mathbf{p}_m^{\text{start}} + \delta \varphi_m$. Each sample contains N = 32 points and K = 12 random samples are generated with the same parameters.

3.1 Discussion

Geochemical consensus scores were computed for the validation data using the algorithm described in Sec. 2.2.3. In Fig. 9, $\hat{f}_{\text{consensus}}$ is plotted as a function of η (the fraction of lower cluster points in each sample) and each curve is parameterized by δ (the magnitude of the outliers which is proportional to the distance between the sampling ellipses). These results are reassuring as the curves generally exhibit monotonic decreasing behaviour with respect to the distortion δ and amount of outliers as η (or $1 - \eta$) increases. This is encouraging since neither δ nor η is known in practice. However, this trend is disrupted in the region $\eta \in [0.4, 0.6]$ due to a breakdown in the MCD robust estimator [9]. It becomes ineffective when η exceeds 0.4 and produces unreliable (ultra-conservative) estimates of $\hat{f}_{\text{outliers}}$ and $d_{\text{outliers}}^{\text{gmean 4}}$. This in turn means the geochemical consensus score, as it is currently formulated, underestimates the impact of outliers in situations where there is a roughly even mix of inliers and outliers. A remedy is proposed in Sec. 4 to

⁴The 0.4 figure is dependent on the sample size N. The MCD robust estimator can handle at most $(N - (n_{\text{variables}} + 1))/2$ outliers before breakdown occurs. When N = 32 and $n_{\text{variables}} = 2$, this evaluates to $\lfloor 14.5 \rfloor$ which translates to 0.4375 of N.



Figure 9: Validation shows monotonic behaviour over the range $\eta \in [0, 0.4) \cup (0.6, 1]$

overcome this limitation. This essentially involves applying the MCD estimator to split-sequences and diminishing sets. With this modification, the results in Fig. 10 show the geochemical consensus measure responds appropriately in the $\eta \in [0.4, 0.6]$ range. This allows the impact of outliers to be quantified when the data contains an unknown and possibly a high level of contamination, especially when outliers account for more than 40% of a sample.



Figure 10: With modification, validation shows reasonable behaviour when contamination-level $\eta \ge 0.4$

4 Extension: split-sequence analysis to address estimator breakdown

When the configuration option handle_breakdown is enabled, the following procedure is performed after η_{outliers} , $d_{\text{outliers}}^{\text{gmean}}$, $d_{\text{outliers}}^{\text{unmasked}}$ and $f_{\text{consensus}}^{\text{existing}}$ are computed in equations (10)–(13). The main ideas behind this algorithm are explained in Section 4.1, then an example is provided in Section 4.2 to illustrate how this works.

Algorithm Split-sequence analysis

Require: Sample where points can be sorted according to a dominant trend **Assume:** Indices start from 1 (not 0). **Input:** Chemical composition $\mathbf{c} \in \mathbb{R}^{N \times m}_+$, ILR-transformed values $ilr(\mathbf{c}) \in \mathbb{R}^{N \times (m-1)}$. Existing estimates: $\hat{d}_{\text{outliers}}^{\text{gmean}}$, $\hat{d}_{\text{outliers}}^{\text{unmasked}}$, $\hat{\eta}_{\text{outliers}}$, $\hat{f}_{\text{consensus}}^{\text{existing}}$ Tolerance: τ (Default: 0.5) Greedy search: *terminate_early* (Default: True) 1: Sort c by dominant element in ascending order π (Default: use c[:, 1] for sorting, PCA projection is also viable) 2: Reorder rows: $\mathbf{x} \leftarrow \pi(\mathbf{c}), \mathbf{y} \leftarrow \pi(\operatorname{ilr}(\mathbf{c}))$ 3: Compute $\mu = \text{median}(\mathbf{x}[:,1])$ and $d_{\text{ref}} = \max\{\hat{d}_{\text{outliers}}^{\text{gmean}}, \sqrt{\chi_{p,\nu}^2}\}$ 4: Set offset $\alpha = 0$ and *split_outliers_assimilate_well* = False while $\alpha \leq 1$ and *split_outliers_assimilate_well* is False do 5: 6: Compute split-point: $s = (N \gg 1) + \alpha$ Form split-sets: $\mathbf{s}_{\mathbf{x}}^{0} = \mathbf{x}[1:s,:], \mathbf{s}_{\mathbf{x}}^{1} = \mathbf{x}[s+1:N,:].$ Similarly, $\mathbf{s}_{\mathbf{y}}^{0} = \mathbf{y}[1:s,:], \mathbf{s}_{\mathbf{y}}^{1} = \mathbf{y}[s+1:N,:].$ 7: for Split-sequence s in $\{0, 1\}$ do 8: Apply MCD to $\mathbf{s}_{\mathbf{v}}^s \in \mathbb{R}^L$ to obtain robust distance vector $\mathbf{d} \in \mathbb{R}^L$, outlier index vector $\mathbf{i}_{\text{outliers}} \in \mathbb{Z}^{n_{\text{outliers}}}$ 9: 10: if $n_{\text{outliers}} = 0$ then 11: continue end if 12: $d_{\max} = \max(\mathbf{d}[\mathbf{i}_{\text{outliers}}]), \mathbf{i}_{\text{inliers}} = [i \notin \mathbf{i}_{\text{outliers}} | 1 \le i \le n_{\text{outliers}}]^T$ 13: if $d_{\max} \ge (2-\tau)d_{\text{ref}}$ then 14: $d_0 = |\mathbf{w}_{average}(\mathbf{s}_{\mathbf{x}}^s[\mathbf{i}_{outliers}, 1], \mathbf{d}[\mathbf{i}_{outliers}]) - \mu|, d_1 = |\mathbf{w}_{average}(\mathbf{s}_{\mathbf{x}}^s[\mathbf{i}_{inliers}, 1], \mathbf{d}[\mathbf{i}_{inliers}]) - \mu|$ 15: 16: if $d_0 < d_1$ then Apply MCD to vstack $(\mathbf{s}_{\mathbf{v}}^{1-s}, \mathbf{s}_{\mathbf{v}}^{s}[\mathbf{i}_{outliers}])$ to obtain robust distance vector $\mathbf{d} \in \mathbb{R}^{M}$ 17: 18: Compute evidence strength: $\xi = \log_{10} \left((d_{\text{max}}/d_{\text{ref}})^{1.5} (d_1/d_0) \right)$ 19: if $\xi < 0.3$ and $d_{\rm max} < 2.5 \times d_{\rm ref}$ then 20: continue end if 21: Compute relaxation constant: $\lambda = \max\{1, \min\{\xi, 2.5\}\}$ 22: if $\min(\mathbf{d}[M - n_{\text{outliers}} + 1 : M]) < \max\{\lambda \cdot d_{\text{ref}}, \frac{1}{2}d_{\max}\}$ then 23: *split* outliers assimilate well = True 24: 25: break 26: end if 27: end if end if 28: 29: end for 30: $\alpha \leftarrow \alpha + 1$ 31: end while (to be continued on page 14)

4.1 Intuition

The algorithm comprises two parts. The first part (line 1–31) reorders points within the sample and performs MCD on half-sequences to avoid misdetection and breakdown in the robust estimator when the number of outliers exceed the $\frac{1}{2}(N-n_{\text{variables}}-1)$ limit [nominal ratio of $\eta \sim 40\%$]. Fig. 11 illustrates a situation where the blue points (9/16) are drawn from cluster 1 and the gray points (7/16) are drawn from cluster 2. Fig. 11(a) describes the reordering step in lines 1–2 of the Algorithm, where the points are sorted by the dominant component (e.g. concentration of Fe or first PCA coefficient) based on prior analysis. With an outlier ratio of 45%, the MCD estimator would likely fail in this case. However, if the ordered sequence is split into halves as shown in Fig. 11(c) and the MCD is applied separately to s_v^0 (left) and s_v^1 (right), the incongruous point circled in Fig. 11(d) may be successfully detected. In line 6, an offset is



Figure 11: Motivation for part 1 of the algorithm

added to the splitting point to handle the special case where the blue/gray point ratio is 50:50. Subsequently, the sample is further considered only if two conditions are satisfied. The first, stipulated in line 14, requires the incongruous point (i_{outliers}) discovered in $\mathbf{s}_{\mathbf{y}}^s$ to be a significant outlier. The second, specified in lines 15–16, requires the composition of $\mathbf{c}[i_{\text{outliers}}, 1]$ to be closer to the median than inliers within the split $\mathbf{s}_{\mathbf{y}}^s$. The second condition helps eliminate cases where $\eta < 25\%$ which warrants no special treatment. In lines 23–24, MCD is applied to $\{\mathbf{s}_{\mathbf{y}}^{1-s}, \mathbf{s}_{\mathbf{y}}^s[i_{\text{outliers}}]\}$, and a decision is made on whether i_{outliers} assimilates well with the other half, $\mathbf{s}_{\mathbf{y}}^{1-s}$. Samples proceed to part 2 of the assessment upon receiving a positive decision.

The second part (line 32–64) applies MCD to an augmented half-sequence, $\{\mathbf{s}_{\mathbf{y}}^{1-s}, \mathbf{s}_{\mathbf{y}}^{s}[\mathbf{r}[1:L-n_{\text{evict}}]]\}$. This process involves eliminating progressively more and more elements, until an improved estimate of $f_*^{\text{consensus}}$ is found or the maximum number of evictions, $n_{\text{evictions}}^{\max}$, is reached. Fig. 12(e) shows the elimination order is established by vector $\mathbf{r} \in \mathbb{R}^L$ which ranks the elements in $\mathbf{s}_{\mathbf{y}}^{s}$ (e.g. right half), excluding i_{outliers} , based on their disparity with respect to $\mathbf{s}_{\mathbf{y}}^{1-s}$ (e.g. left half). Sorting in line 38 uses the robust distances d computed for points in $\mathbf{s}_{\mathbf{y}}^{s}$. Fig. 12(f) shows the interim estimates obtained from diminishing sets. In earlier iterations, the MCD estimator may struggle to find outliers or report a meaningful $d_{\text{outliers}}^{\text{gmean}(k)}$ value larger than the existing estimate. With successive eviction, the proportion of outliers is diminished and the situation steers further away from breakdown. The final fraction of outliers, $\eta_{\text{outliers}}^{(\text{final})}$, includes both the number of evicted elements ($n_{\text{evict}} \equiv k$) and outliers ($n_{\text{outliers}}^{(k)}$) reported in the reduced set $\{\mathbf{s}_{\mathbf{y}}^{1-s}, \mathbf{s}_{\mathbf{y}}^{s}[\mathbf{r}[1:L-k]]\}$.



Figure 12: Sketch for part 2 of the algorithm

Algorithm Split-sequence analysis (continued)

32: Apply MCD estimator to diminishing sets 33: if *split_outliers_assimilate_well* is True then 34: Set $\mathbf{d}[\mathbf{i}_{\text{outliers}}] = 0$ Set $n_{\text{iterations}}^{\text{max}} = 8$, $n_{\text{evictions}}^{\text{max}} = N \gg 2$, $\Delta = \max\{1, \lfloor n_{\text{evictions}}^{\text{max}}/n_{\text{iterations}}^{\text{max}} \rfloor\}$. 35: Initialise: $i = 0, n_{\text{evict}} = 1 - \Delta, \hat{f}_*^{\text{consensus}} = []$ 36: Set $\hat{f}_{\text{consensus}} = \hat{f}_{\text{consensus}}^{\text{existing}}$ 37: Compute retention vector: $\mathbf{r} = \operatorname{argsort}(\mathbf{d}, \operatorname{ascending}) \in \mathbb{Z}^{L}$ 38: while $n_{\text{evict}} < n_{\text{evictions}}^{\max}$ do 39: 40: $n_{\text{evict}} \leftarrow n_{\text{evict}} + \Delta$ Apply MCD to vstack $(\mathbf{s}_{\mathbf{y}}^{1-s}, \mathbf{s}_{\mathbf{y}}^{s}[\mathbf{r}[1:L-n_{\text{evict}}]])$ to obtain 41: robust distance vector $\tilde{\mathbf{d}} \in \mathbb{R}^{N-n_{\text{evict}}}$ and outlier index vector $\tilde{\mathbf{i}}_{\text{outliers}} \in \mathbb{Z}^{n_{\text{outliers}}}$. 42: if $n_{\text{outliers}} = 0$ then continue 43: 44: end if Compute $d_{\text{robust}} = \exp\left(\frac{1}{n_{\text{outliers}}}\sum_{i \in \tilde{\mathbf{i}}_{\text{outliers}}}\log\tilde{\mathbf{d}}[i]\right)$ 45: if $d_{\text{robust}} < \sqrt{\chi^2_{p,\nu}}$ then 46: continue 47: end if 48: $i \leftarrow i+1$ Set $d_*^{\text{gmean}}[i] = d_{\text{robust}}$ 49: 50: Set $d_*^{\text{unmasked}}[i] = \min\{\log_{10}\left(d_*^{\text{gmean}}[i]/\sqrt{\chi_{p,\nu}^2}\right), 1\}$ 51: Set $\eta_*^{\text{outliers}}[i] = \min\{f, 1 - f\}$ where $f = (n_{\text{outliers}} + n_{\text{evict}})/N$ Set $f_*^{\text{consensus}}[i] = (1 - \eta_*^{\text{outliers}}[i])^{d_*^{\text{unmasked}}[i]}$ 52: 53: if $terminate_early$ and $f_*^{\text{consensus}}[i] < 0.8 \times f_{\text{consensus}}^{\text{existing}}$ then 54: break 55: end if 56: end while 57: if $f_*^{\text{consensus}}$ is not empty then 58: $k = \arg\min_i f_*^{\text{consensus}}$ 59: $\begin{aligned} & \text{if } f_*^{\text{consensus}}[k] < \hat{f}_{\text{consensus}}^{\text{existing}} \text{ then} \\ & \text{Set } \hat{d}_{\text{outliers}}^{\text{gmean}} \leftarrow d_*^{\text{gmean}}[k], \hat{d}_{\text{outliers}}^{\text{unmasked}} \leftarrow d_*^{\text{unmasked}}[k], \hat{\eta}_{\text{outliers}} \leftarrow \eta_*^{\text{outliers}}[k], \hat{f}_{\text{consensus}} \leftarrow f_*^{\text{consensus}}[k] \end{aligned}$ 60: 61: 62: end if end if 63: 64: end if **Output:** $\hat{d}_{\text{outliers}}^{\text{gmean}}, \hat{d}_{\text{outliers}}^{\text{unmasked}}, \hat{\eta}_{\text{outliers}}, \hat{f}_{\text{consensus}}$

4.2 Example

This example is based on the test data ("high-breakdown.csv") supplied with the source code. The input and relevant quantities shown on lines 1–3 of the algorithm are specified in Table 3.

| Sample contains $N = 16$ points, seven of which (43.75%) are outliers. With $n_{\text{variables}}(\mathbf{y}) = 2$, at most $\lfloor (N - n_{\text{variables}} - 1)/2 \rfloor = 6$ outliers can be detected | | | | | | | | | | | | | | etected. | | |
|--|-------|-------|--------|-------|-------|-------|-------|-------|-------|-------|-----------|-----------|-------|-----------|-------|-----------|
| Point | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | <u>11</u> | <u>12</u> | 13 | <u>14</u> | 15 | <u>16</u> |
| $Fe \leftrightarrow c[:, 1]$ | 0.912 | 0.861 | 0.863 | 0.891 | 0.914 | 0.883 | 0.894 | 0.873 | 0.889 | 0.430 | 0.441 | 0.468 | 0.429 | 0.468 | 0.481 | 0.426 |
| $SiO_2 \leftrightarrow c[:, 2]$ | 0.058 | 0.079 | 0.073 | 0.048 | 0.049 | 0.079 | 0.052 | 0.072 | 0.060 | 0.299 | 0.323 | 0.296 | 0.303 | 0.312 | 0.271 | 0.307 |
| $Al_2O_3 \leftrightarrow \mathbf{c}[:,3]$ | 0.030 | 0.060 | 0.064 | 0.062 | 0.037 | 0.038 | 0.054 | 0.055 | 0.051 | 0.271 | 0.236 | 0.237 | 0.268 | 0.220 | 0.248 | 0.266 |
| ilr(c) 1 st coeff | 1.950 | 1.690 | 1.751 | 2.069 | 2.070 | 1.710 | 2.009 | 1.765 | 1.905 | 0.257 | 0.219 | 0.325 | 0.247 | 0.286 | 0.406 | 0.231 |
| ilr(c) 2nd coeff | 1.662 | 1.193 | 1.110 | 0.987 | 1.428 | 1.583 | 1.137 | 1.241 | 1.231 | 0.230 | 0.384 | 0.370 | 0.241 | 0.451 | 0.306 | 0.250 |
| Permutation π | 16 | 13 | 10 | 11 | 14 | 12 | 15 | 2 | 3 | 8 | 6 | 9 | 4 | 7 | 1 | 5 |
| $\pi(\mathbf{c}) \rightarrow \mathbf{x}_1$ | 0.426 | 0.429 | 0.430 | 0.441 | 0.468 | 0.468 | 0.481 | 0.861 | 0.863 | 0.873 | 0.883 | 0.889 | 0.891 | 0.894 | 0.912 | 0.914 |
| $\pi(\operatorname{ilr}(\mathbf{c})) \to \mathbf{y}_1$ | 0.231 | 0.247 | 0.257 | 0.219 | 0.286 | 0.325 | 0.406 | 1.690 | 1.751 | 1.765 | 1.710 | 1.905 | 2.069 | 2.009 | 1.950 | 2.070 |
| $\pi(\operatorname{ilr}(\mathbf{c})) \to \mathbf{y}_2$ | 0.250 | 0.241 | 0.230 | 0.384 | 0.451 | 0.370 | 0.306 | 1.193 | 1.110 | 1.241 | 1.583 | 1.231 | 0.987 | 1.137 | 1.662 | 1.428 |
| | | т | 1-1-2. | Tues | | . T | | £ | 1 | | - 1 | | | | | |

 Table 3: Trace part 1: Input, transformed and permuted sequences

Applying MCD as described in Sec. 2.2.3 produces $\hat{\eta}_{\text{outliers}} = 0.25$, $\hat{d}_{\text{outliers}}^{\text{gmean}} = 5.608$, $\hat{d}_{\text{outliers}}^{\text{unmasked}} = 0.3149$ and $\hat{f}_{\text{consensus}} = 0.9133$. The fraction of outliers ($\hat{\eta}_{\text{outliers}}$) and their impact ($\hat{d}_{\text{outliers}}^{\text{gmean}}$) are both underestimated. Table 4 shows the tests performed on the left half sequence, s_y^0 , and the requirements satisfied during the first iteration in the split-sequence FOR loop (see Algorithm: Split-sequence analysis on page 12).

| split-sequences | | | | | s_v^0 | | | $\mathbf{s}_{\mathbf{v}}^{1}$ | | | | | | | | | | |
|---|--------------------|-----------|---------------------|-----------|---------------------------|--------------------------------|---------------|-------------------------------|-------------------|--------------|-----------|---------------|-----------------------|-----------|-------|-------|--|--|
| with offset $\alpha = 0$ | 0.231 | 0.247 | 0.257 | 0.219 | Ŏ.286 | 0.325 | 0.406 | 1.690 | 1.751 | 1.765 | 1.710 | 1.905 | 2.069 | 2.009 | 1.950 | 2.070 | | |
| | 0.250 | 0.241 | 0.230 | 0.384 | 0.451 | 0.370 | 0.306 | 1.193 | 1.110 | 1.241 | 1.583 | 1.231 | 0.987 | 1.137 | 1.662 | 1.428 | | |
| For $s = 0$, | | | | | | | | | | | | | | | | | | |
| robust dist. d | 1.007 | 0.942 | 1.137 | 1.771 | 1.539 | 1.817 | 4.550* | 40.722* | * | | | | | | | | | |
| $ \mathbf{s}_{\mathbf{x}}^{0}[:,1] - \mu_{\mathbf{x}} $ | 0.436 | 0.433 | 0.432 | 0.421 | 0.394 | 0.394 | 0.381 | 0.001 | note: n | nedian μ | x = 0.86 | 2 from > | x ₁ | | | | | |
| i _{outliers} (*) | [7, 8] | | | | | | | | | | | | | | | | | |
| statistics | $d_{\text{max}} =$ | 40.722 | , $d_{\rm ref} = 5$ | 6.608 and | $1 \tau = 0.5$ | 6. Requir | ement 1 [| $d_{\max} \ge (2$ | $-\tau d_{ref}$] | is satisf | ied. | | | | | | | |
| | weight | ed distan | ces from | median | for \mathbf{i}_{outlie} | $_{rs}$ and \mathbf{i}_{inl} | iers: $d_0 =$ | $0.039, d_1 =$ | = 0.414. | Require | ment 2 [a | $l_0 < d_1$] | is also s | atisfied. | | | | |

 Table 4: Trace part 2: Tests performed on left half sequences

Applying MCD to the augmented half-sequence, $\mathbf{s}_{\mathbf{y}}^{\dagger} \stackrel{\text{def}}{=} (\mathbf{s}_{\mathbf{y}}^{0}[\mathbf{i}_{\text{outliers}}], \mathbf{s}_{\mathbf{y}}^{1})$ as per Algorithm (line 17), Table 5 shows the values obtained. The incongruous point assimilates well with $\mathbf{s}_{\mathbf{y}}^{0}$; thus it qualifies for part two of the assessment (lines 32–64). Table 6 shows the elimination order, the formation of diminishing sets, and the values computed after each element is evicted.

| · • • | | 0.400 | 1 (00 | 1 7 7 1 | 1 7 6 7 | 1 7 1 0 | 1.007 | 2.0(0 | 2 000 | 1.050 | 0.070 | | |
|--|--|--------|---------|---------|---------|---------|-------|-------|-------|-------|-------|--|--|
| augmented seq s_v^{\dagger} | | 0.406 | 1.690 | 1./51 | 1.765 | 1./10 | 1.905 | 2.069 | 2.009 | 1.950 | 2.070 | | |
| | | 0.306 | 1.193 | 1.110 | 1.241 | 1.583 | 1.231 | 0.987 | 1.137 | 1.662 | 1.428 | | |
| robust dist. $\tilde{\mathbf{d}}$ | | 11.667 | *1.434* | 1.281 | 0.847 | 1.738 | 0.299 | 1.829 | 1.078 | 1.896 | 1.541 | | |
| evidence strength | $\xi = \lambda = 2.313$ | | | | | | | | | | | | |
| | The condition $\min(\tilde{\mathbf{d}}[1:n_{\text{outliers}}]) < \max\{\lambda \cdot d_{\text{ref}}, \frac{1}{2}d_{\max}\}$ is satisfied. | | | | | | | | | | | | |
| | $split_outliers_assimilate_well$ is set to True. The tests for $s = 1$, and $\alpha = 1$ (with $s \in \{0, 1\}$) will not be performed. | | | | | | | | | | | | |
| Table 5: Trace part 3: Applying MCD to augmented half-sequence | | | | | | | | | | | | | |

| robust dist. d' | 1.007 | 0.942 | 1.137 | 1.771 | 1.539 | 1.817 | 0* | 0* | | | | | | | | |
|---|--|-----------------|------------------|---------|----------|---------|----------|---------|--------------|-------|--------|----------|---------|--------------|--------------|----|
| priority r | 7 | 8 | 2 | 1 | 3 | 5 | 4 | 6 | | | | | | | | |
| diminishing sets | $(\mathbf{s}_{\mathbf{y}}^{1}, \mathbf{s}_{\mathbf{y}}^{0}[\mathbf{r}[1:L-n_{\text{evict}}]])$. For clarity, elements retained from $\mathbf{s}_{\mathbf{y}}^{0}$ are underlined. | | | | | | | | | | | | | | | |
| $(n_{\text{evict}} = 1)$ | 1.751 | 1.765 | 1.710 | 1.905 | 2.069 | 2.009 | 1.950 | 2.070 | <u>0.406</u> | 1.690 | 0.247 | 0.231 | 0.257 | 0.286 | <u>0.219</u> | × |
| | 1.110 | 1.241 | 1.583 | 1.231 | 0.987 | 1.137 | 1.662 | 1.428 | 0.306 | 1.193 | 0.241 | 0.250 | 0.230 | <u>0.451</u> | <u>0.384</u> | × |
| robust dist. $\tilde{\mathbf{d}}$ | 1.281 | 0.847 | 1.738 | 0.299 | 1.829 | 1.078 | 1.896 | 1.541 | 11.667* | 1.434 | 12.835 | 5*12.919 | *12.793 | *12.137 | *12.706 | 5* |
| <i>i</i> _{outliers} | [9,11, | 12,13,14 | 4,15] | | | | | | | | | | | | | |
| adjusted f | (noutlie | $ers + n_{ers}$ | $_{\rm vict})/N$ | =(6+ | 1)/16 = | = 0.437 | 5 | | | | | | | | | |
| $\hat{\eta}_{\text{outliers}}$ | min{; | f, 1 - f | * = 0.4 | 375 | | | | | | | | | | | | |
| $\hat{d}^{\text{gmean}}_{\text{outliers}}$ | 12.50 | 0932 | - | | | | | | | | | | | | | |
| $\hat{d}_{\text{outliers}}^{\text{unmasked}}$ | 0.6629 | 980 | | | | | | | | | | | | | | |
| $\hat{f}_{consensus}$ | 0.682 | 867 | Further | evictio | n is not | require | d as a s | olution | is found. | | | | | | | |

Table 6: Trace part 4: Elimination order, formation of diminishing sets and the output after one element is evicted

Here, the process terminates after one eviction. The estimated fraction of outliers $\hat{\eta}_{\text{outliers}} = 0.4375$ corresponds exactly to the ground truth. It also produces a reasonable geochemical consensus value of 0.6828; as opposed to 0.9133 when the estimator breaks down—which is misleading. This example and the graphical results presented in Fig. 10 reaffirm the viability of the consensus score and its ability to measure outlier distortion with the split-sequence extension when the contamination level is unknown and possibly as high as 50%.

5 Implementation

A Python implementation of the proposed algorithms will be available on GitHub at https://github.com/raymondleung8/sampling-consensus. The code will be distributed open-source under a BSD 3-clause license.

6 Conclusion

The mining industry currently lacks a standard framework that can convey the degree of confidence in the average grade estimates computed for selective mining units (otherwise known as grade-blocks). This paper proposes two measures for quantifying sampling fairness and geochemical consensus for blasthole assays located within a grade-block. Sampling fairness considers spatial factors such as the sampling density and bias in the spatial distribution of the blastholes. Geochemical consensus considers the disparity between the assay samples collected from a grade-block. A novel masking expression (12) is used to measure meaningful differences. This takes into account the fraction of outliers observed (10) and magnitude of distortion above a threshold (11). An extension based on split-sequence analysis has been proposed to address the issue of MCD estimator breakdown when the fraction of outliers exceeds $(n_{\text{samples}} - n_{\text{variables}} - 1)/(2n_{\text{samples}})$. This allows the true extent and impact of outliers to be quantified, and the consensus score to remain effective when samples contain up to 50% of outliers. The method was tested on data obtained from a banded iron formation (BIF) hosted iron ore deposit within the Hamersley region in Western Australia. An open-source implementation of the algorithms is available on github.

Acknowledgements

This work was supported by the Australian Centre for Field Robotics and the Rio Tinto Centre for Mine Automation.

References

- [1] Piotr Szmigiel. *Discovery of an optimum selective mining unit for a high quality selective deposit.* PhD thesis, University of Missouri-Rolla, 2005.
- [2] SG Hagemann, T Angerer, Paul Duuring, CA Rosière, RC Figueiredo e Silva, L Lobato, AS Hensler, and DHG Walde. BIF-hosted iron mineral system: a review. *Ore Geology Reviews*, 76:317–359, 2016.
- [3] Raymond Leung, Alexander Lowe, Anna Chlingaryan, Arman Melkumyan, and John Zigman. Bayesian surface warping approach for rectifying geological boundaries using displacement likelihood and evidence from geochemical assays. Available at: https://doi.org/10.1145/3476979. ACM Transactions on Spatial Algorithms and Systems, 8(1):1–23, March 2022.
- [4] R Hargreaves, G.W. Elkington, T. Booth, and W.J. Shaw. Mine reconciliation standardisation R factor series. In *International Mining Geology Conference*, pages 366–374. AusIMM, 2022.
- [5] Mia Hubert and Michiel Debruyne. Minimum covariance determinant. *Wiley interdisciplinary reviews: Computational statistics*, 2(1):36–43, 2010.
- [6] Raimon Tolosana-Delgado, Ute Mueller, and K Gerald van den Boogaart. Geostatistics for compositional data: an overview. *Mathematical Geosciences*, 51(4):485–526, 2019.
- [7] Michail T Tsagris, Simon Preston, and Andrew TA Wood. A data-based power transformation for compositional data. In J Egozcue, R Tolosana-Delgado, and M Ortego, editors, *Proceedings of the 4th International Workshop* on Compositional Data Analysis, pages 565–572. Springer, 2011.
- [8] HO Lancaster. The Helmert matrices. The American Mathematical Monthly, 72(1):4-12, 1965.
- [9] Mia Hubert, Peter J Rousseeuw, and Stefan Van Aelst. High-breakdown robust multivariate methods. *Statistical Science*, 23(1):92–119, 2008.