# Author's Accepted Manuscript

Variable lag variography using k-means clustering

I.K. Kapageridis

Cite this article as: I.K. Kapageridis, Variable lag variography using k-means c l u s t e r i n g , *Computers and Geosciences,* http://dx.doi.org/10.1016/j.cageo.2015.04.004

# Variable Lag Variography Using k-means Clustering

I.K. Kapageridis

Laboratory of Mining Information Technology and GIS Applications, Department of Environmental Engineering, Technological Educational Institute of Western Macedonia, Greece, Email: ioannis.kapageridis@gmail.com

## Abstract

Experimental variography in three dimensions based on drillhole data and current modelling software requires the selection of particular directions (azimuth and plunge) and a basic lag distance. Variogram points are then calculated on distances which are multiples of that basic lag. As samples rarely follow a regular grid, directional and distance tolerances are applied in order to have sufficient pairs to calculate reliable variogram points. This process is adequate when drillholes follow a drilling pattern (even if not an exactly regular grid) but can be time consuming and hard when the drilling pattern is irregular or when drillhole orientations vary considerably. Having all variogram points being calculated on multiples of a fixed lag, and the same tolerance being applied throughout the range of distances used, can be very restrictive and a reason for considerable time wasting or even failure to calculate an interpretable experimental variogram. The method discussed in this paper is using k-means clustering of sample pairs based on pair separation distance leading to a number of clusters each representing a different variogram point. This way, lag parameters are adjusted automatically to match the spatial distribution of sample locations and the resulting variogram is improved. Case studies are provided showing the benefits of this method over current fixed-lag experimental variogram calculation techniques.

keywords: experimental variogram, k-means clustering, variogram modelling

## Introduction

Drilling patterns, in mineral exploration programmes in particular, very rarely follow a strictly regular pattern. In some cases, this is due to practical issues causing a deviation from a constantly spaced pattern, while in other cases, it is the geometry of the targeted orebody envelope that requires a more flexible pattern to be followed. As exploration takes place in stages, in-fill drilling to increase the level of confidence in certain areas, also causes local changes in sample spacing. Drilling from underground workings is, in most cases, irregular and leads to extreme variation of sample spacing.

Regardless of the reason and the degree of irregularity, when the sample spacing is not reasonably constant, the variography practitioner can face a very difficult task in finding a set of lag parameters that work for all variogram points in a particular direction. In most cases, parameters that work for the variogram points at smaller separation distances will not

37  work for the points at larger separation distances and vice versa. The problem is further
38  exaggerated by the way the lag parameters are applied by geostatistical software.
39  Depending on how dynamic is the application of the parameters, it is possible to reach, after
40  a considerable amount of time and effort, a set of parameters that works reasonably well for
41  most variogram points. However, some of the geostatistical software is not dynamic at all in
42  the application of lag and direction parameters. Setting of these parameters and running the
43  computation of experimental variogram values are independent and take place separately,
44  in which case, it can be almost impossible to find a set of parameters that will produce an
45  interpretable experimental variogram.

46      The problem of grouping pairs of samples that fall in a particular direction seems to
47  be quite suitable for solving using a clustering algorithm like the k-means. The main concept
48  is that pairs are selected to belong to a particular direction according to the vector they
49  define and then they are grouped according to a criterion like the separation distance. This
50  way, the user does not have to spend time and effort in finding appropriate lags and
51  tolerances that work in that particular direction. A separate clustering run will have to be
52  performed for each direction considered. The directions themselves could be chosen using a
53  similar approach to better match the most sampled directions, but this is something that
54  most geostatistics practitioners would probably like to keep control of, and computationally
55  it would require a lot more time to perform.

# Current State of Variogram Calculation

57  The information provided in this section is based on the experience and knowledge of the
58  author and does not necessarily cover every geostatistical software package and method
59  available to the geostatistics practitioner. However, the author believes that most of the
60  available packages follow, in some way, one of the paradigms discussed here.

## Variogram Points Based on Multiples of a Basic Lag and an Absolute Lag Tolerance
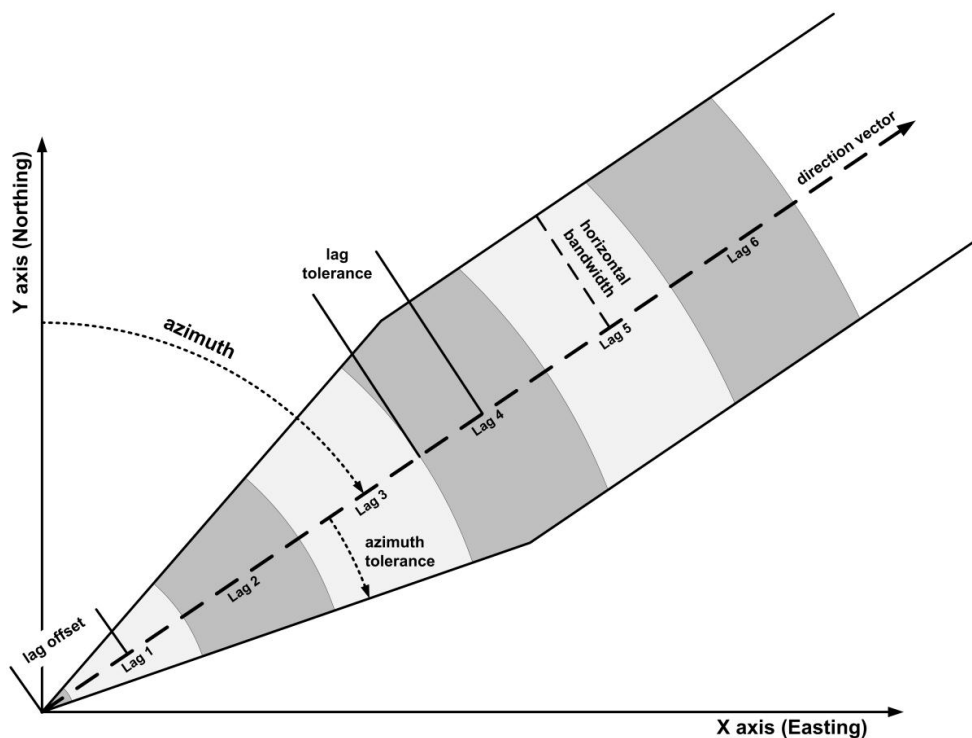
63  Most geostatistical software packages follow this concept. As shown in Figure 1 (for a 2D
64  case), for a particular direction chosen, a search area is defined using some direction
65  tolerance which can be controlled both horizontally and vertically. Some packages use the
66  same direction tolerance in both cases while others allow for separate tolerances to be
67  applied for azimuth and plunge. These tolerances are allowed to expand the search area as
68  the separation distance increases up to a maximum distance (bandwidth) from the direction
69  vector. This way, the search area begins as a cone of circular or elliptical section (depending
70  on whether the azimuth and plunge tolerances are different), and then becomes a cylinder
71  of similar section to the cone, once the maximum distance from the direction vector is
72  reached. Some packages allow for a separate maximum distance to be applied horizontally
73  and vertically.

74      The search area is split into multiple search windows which are defined using a basic
75  lag and a lag tolerance. Each search window is centred on a distance along the direction
76  vector derived by a multiple of the basic lag plus some lag offset. The extents of the search

77  window are controlled by the lag tolerance, which, in most cases, is fixed to an absolute
78  distance value and does not change with distance. For example, if the basic lag is 50m and
79  the tolerance is 15m, then the search window at the sixth variogram point will be centred at
80  6x50 = 300m and start at 285m and stop at 315m (for a zero lag offset). The sample pairs
81  that fall within the search area are checked against the search windows and they get
82  grouped into different variogram points according to the corresponding separation distance.

83  As a percentage of the separation distance, the tolerance decreases with every
84  multiple of the basic lag, leading to an ever decreasing number of pairs found at higher
85  distances. In the previous example, for the first variogram point, 15m is 30% of the 50m
86  window distance, 15% of the 100m, 10% of the 150m, and so on. Of course, this is not the
87  only reason for the number of pairs to decrease with distance - it will happen inevitably as
88  we reach the maximum distance between drillholes. However, it is probably the only reason
89  or factor that can be controlled using a different approach of searching for pairs, i.e. a
90  different way of defining the search windows. It should be noted that for variogram
91  smoothing purposes, some of the packages allow overlapping of the search windows, in
92  which case, some sample pairs are used in more than one variogram points.

93  Some of the main examples of geostatistical packages (the list is far from comlete)
94  that follow the approach described above in two or three dimensions are Isatis (Bleines et al,
95  2004), GSLIB (Deutsch et al, 1992) including the implementation for standard directional
96  variography in Vulcan 3D software (Maptek Pty Ltd), GEO-EAS (Englund et al, 1991), SGeMS
97  (Remy et al, 2011), and VarioWin (Pannatier, 1996). In some of these software packages, a
98  file containing all pairs of samples and their separation vector parameters (distance,
99  azimuth, plunge) is formed before the variogram points are calculated (called a pair
100 comparison file).



101

102  Figure 1: Standard sample pair selection in two dimensions, used in most current
103  geostatistical software.

## Variogram Points Based on 3D Block Search Windows

105  This method is based on orthogonal blocks forming a 3D model centred on the origin of
106  variography polar coordinates – the model always has an odd number of blocks along X, Y
107  and Z. The centroid of each block together with the origin defines a different vector with its
108  own azimuth and plunge. The block extents control the lag, azimuth and plunge tolerance
109  but in a way quite different to the technique described in the previous section. Each of the
110  sample pairs is checked against each block, with one sample at the model origin. If the
111  second sample of the pair falls is nearer the centroid of a particular block, the pair is used to
112  calculate the variogram value for that block (Figure 2). Once all pairs are assigned to
113  particular blocks, the block variogram values are calculated and stored in the model. The end
114  result is a three-dimensional variogram map that can be displayed using a number of
115  different methods (contours, slices, shells, etc.) in two or three dimensions. In addition to
116  the variogram value (different variogram types are available such as standard
117  semivariogram, general relative, pairwise relative, etc.), some other useful information is
118  calculated and store in each block, including the number of pairs, average distance, average
119  head and tail values of the pairs.



120

121 Figure 2: Simple 2D representation of the way sample pairs are selected for the calculation
122 of a particular block variogram value in 3D block variography. Blocks are coloured by
123 variogram value and can potentially reveal the orientation of anisotropy.

124       The quality of the produced variogram map is controlled by the block sizes and
125 tolerances along X, Y, and Z relative to the sample spacing. Block sizes work similarly to the
126 lag sizes in the previous technique. The tolerances along X, Y, and Z work relative to the
127 block centroids (blocks can be allowed to overlap when checked against sample pairs). As an
128 approach, it is more suitable to investigating the existence and orientation of geometrical
129 (ellipsoidal) anisotropy, and finding the directions that work better with the available
130 sampling pattern, rather than forming the basis for variogram modelling. This technique,
131 called cube variography, is available in the more recent versions of Vulcan 3D (Maptek Pty
132 Ltd).

133       As a technique, it is still not particularly dynamic in its application, as the block
134 model has to be calculated first using a parameter file that needs to be modified and called
135 again if a different block setup is necessary. However, it is an improvement in the visual
136 aspect of checking the effects of different block sizes (i.e. the effect of different lag sizes) in
137 different directions, potentially leading to a better lag setup in directional variography. The
138 orientation of the anisotropy ellipsoid can be easily determined and the number of
139 directional variograms to calculate can be potentially reduced. At its current form, this
140 technique can work as a preparatory step before the common directional variography
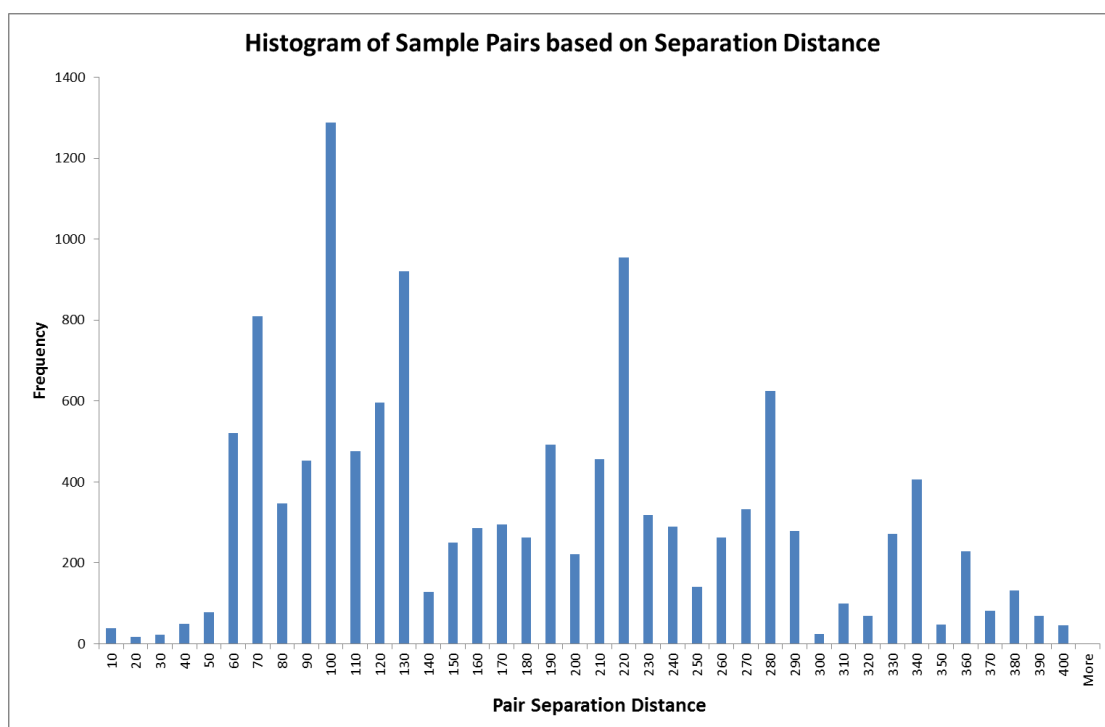141 described in the previous section.

## 142 Variable Lag Experimental Variogram Calculation Based on k-
## 143 means Clustering

### 144 Concept

145 The techniques described in the previous sections require a time consuming trial and error
146 procedure to be followed in order to reach a lag setup that will produce a reasonably
147 interpretable experimental variogram. The information to reach a good lag setup is already
148 available in the sample pairs for any particular direction, even for irregular sampling
149 patterns, but it is probably too much for the practitioner to handle. Figure 3 shows a
150 histogram of sample pairs based on separation distance for the data of the first case study.
151 Some separation distances present much higher pair frequencies than others, and, together
152 with some very low frequency distances, they can define a group of pairs that could produce
153 a reliable experimental variogram point with sufficient number of pairs. For example, such a
154 point could be considered between 250 and 300 meters around the peak at 280. Other
155 points can be identified in a similar manner. The end result would be a set of experimental
156 variogram points defined at variable separation distances, not multiples of a basic lag, and
157 with varying distance tolerances – a concept called variable lag variography (VLV) from this
158 point on (Figure 4). These points would have sufficient pairs to be considered reliable even
159 at higher separation distances.

160    Following such logic in a manual way, by examining a histogram like the one in
161    Figure 3, could be beneficial but would still require a fair amount of time and work. It would
162    be better if these groups of pairs can be identified by an automated procedure, as
163    unsupervised as possible. The author has chosen k-means clustering for its speed and
164    simplicity. It is not necessarily the best method for this problem and still requires some
165    minimum input by the user. There are many similar clustering methods and variations, and it
166    is one of the aims for future work to identify one that is more appropriate and will produce
167    the most optimum results. However, k-means proved sufficient to demonstrate the validity
168    of the VLV concept. IBM SPSS Statistics, the software package used for clustering in this
169    study, provides two more clustering methods, the TwoStep Cluster Analysis, and the
170    Hierarchical Cluster Analysis (IBM SPSS Statistics Base 20, 2011).

171



172

173    Figure 3: Histogram of sample pairs for a particular direction based on separation distance in
174    the case of the underground drilling pattern of case study 1.

175

176    Variography parameters using k-means clustering in VLV are represented by the resulting
177    clustering information:

178    • Lag offset: the average separation of the first cluster (first variogram point).
179    • Lag: the average separation of each cluster (each variogram point) - different for
180      each variogram point, not a multiple of a standard distance.
181    • Lag tolerance: the maximum distance of the pairs classified in a specific cluster from
182      that cluster's center - different for each variogram point, not a fixed value.
183    • Pair count: the number of pairs classified in each cluster.

184 Figure 4 shows how these parameters define search windows in the case of VLV in two
185 dimensions.



187 Figure 4: Proposed variable lag sample pair selection based on k-means clustering of pairs.

## k-means Clustering

189 k-means clustering is a method originally used in signal processing, commonly used for
190 cluster analysis in data mining. k-means clustering groups n observations into k clusters, with
191 each observation assigned to the cluster with the nearest mean. The term "k-means" was
192 introduced by MacQueen in 1967. The standard algorithm was first proposed by Lloyd in
193 1957 as a technique for pulse-code modulation. Forgy published essentially the same
194 method (Forgy, 1965), which is why it is sometimes referred to as Lloyd-Forgy. A more
195 efficient version was proposed and published by Hartigan and Wong in 1975 and 1979.

196 It is an iterative algorithm that is performed in steps. Before any iteration, the
197 clusters are initially centred on an equal number of observations. These observations can be
198 chosen using different methods. Iterations involve two steps. In the first step, each
199 observation is assigned to the cluster whose mean yields the least within-cluster sum of
200 squares. The second step involves the calculation of the new means to be the centroids of
201 the observations in the new clusters. The algorithm converges when there is no change in
202 the assignments.

203 The implementation of the k-means clustering algorithm in SPSS (called QUICK
204 CLUSTER) can handle large numbers of cases. It attempts to identify relatively homogeneous
205 groups of cases based on selected characteristics. As with most k-means clustering

206 algorithms, it requires that the number of clusters is specified a priori. The initial cluster
207 centres can be manually selected if required. There are two methods available for classifying
208 cases, either updating cluster centres iteratively or classifying only. Information such as
209 cluster membership, distance information, and final cluster centres can be stored after
210 clustering. Optionally, a variable can be specified whose values will be used to label case-
211 wise output. The first iteration of the algorithm involves three steps (as described in IBM
212 SPSS Statistics 20 Algorithms, 2011):

213 Step 1 - Initial Cluster Centre Selection
214 Selection of the initial cluster centres involves a single pass of the data. The values of the
215 first NC (number of requested clusters) cases are selected as cluster centres, and the
216 remaining cases are reprocessed as follows:

    217     a) If $\min_i d(x_k, M_i) > d_{mn}$ and $d(x_k, M_m) > d(x_k, M_n)$, then $x_k$ replaces $M_n$. If $\min_i d(x_k, M_i) > d_{mn}$
    218     and $d(x_k, M_m) < d(x_k, M_n)$, the $x_k$ replaces $M_m$; that is, if the distance between $x_k$ and its
    219     closest cluster mean is greater than the distance between the two closest means
    220     ($M_m$ and $M_n$), then $x_k$ replaces either $M_m$ or $M_n$, whichever is closer to $x_k$.
    221     b) If $x_k$ does not replace a cluster mean in (a), a second test is made:
    222     ▪ Let $M_q$ be the closest cluster mean to $x_k$.
    223     ▪ Let $M_p$ be the second closest cluster mean to $x_k$.
    224     ▪ If $d(x_k, M_p) > \min_i d(M_q, M_i)$, then $M_q = x_k$;
    225     That is, if $x_k$ is further from the second closest cluster's centre than
    226     the closest cluster's centre is from any other cluster's centre,
    227     replace the closest cluster's centre with $x_k$.

228 where, NC is the number of requested clusters, $M_i$ the mean of the ith cluster, $x_k$ the vector
229 of the kth observation, $d(x_i, x_j)$ is the Euclidean distance between vectors $x_i$ and $x_j$, and $d_{mn}$ is
230 the distance between the two closest means ($\min_{i,j} d(M_i, M_j)$. After one pass through the
231 data, the initial means of all NC clusters are set.

232 Step 2 – Initial Cluster Centres Updating
233 Starting with the first case, each case in turn is assigned to the nearest cluster, and the
234 cluster means are updated. The initial cluster centre is included in this mean. The updated
235 cluster means are considered as the classification cluster centres.

236 Step 3: Cases Assigning to Nearest Cluster
237 The third pass through the data assigns each case to the nearest cluster, where distance
238 from a cluster is the Euclidean distance between that case and the (updated) classification
239 centres. Final cluster means are then calculated as the average values of clustering variables
240 for cases assigned to each cluster. Final cluster means do not contain classification centres.

241     When the number of iterations is greater than one, the final cluster means in step 3
242 are set to the classification cluster means in the end of step 2, and QUICK CLUSTER repeats
243 step 3 again. The algorithm stops when either the maximum number of iterations is reached
244 or the maximum change of cluster centres in two successive iterations is smaller than ε (the
245 convergence criterion) times the minimum distance among the initial cluster centres.

246 ## Application of k-means Clustering to Sample Pairs Grouping for
247 ## Variography

248 As it was mentioned before, in order to test the proposed methodology, two software
249 packages were used: Vulcan 3D, a mine planning package, and IBM SPSS Statistics, a package
250 used for statistical analysis. Vulcan was used to provide the samples database environment
251 and general variography tools for displaying and modelling. Vulcan's Isis database module
252 provides all the necessary functions for manipulating a drillhole or other sample database,
253 while the Envisage graphical environment provides advanced 3D tools for graphically
254 displaying samples. Vulcan's geostatistical functionality is based on GSLIB (Deutsch et al,
255 1992). IBM SPSS Statistics was used to provide the k-means clustering algorithm.

256 A script written by the author in Perl and utilising Vulcan's Lava Perl modules was
257 developed to take care of all the sample pairs preparation work and calling SPSS for
258 clustering. Each direction is processed separately, i.e. the script works in one direction at a
259 time. Lava modules give access to all Vulcan database and model structures as well as the
260 graphical environment through a Perl script. The script allows the user to select the samples
261 database and required sample location and grade fields, as well as the directional
262 parameters for the searching. Currently, it allows the application of an azimuth and plunge
263 with separate tolerances and a bandwidth that is applied in both (Figure 5).

264


265 Figure 5: Specification panel from the script responsible for generating the pairs file for a
266 particular direction and direction tolerances and running SPSS for clustering with k-means.

267

268     The script goes through the following steps when it runs:

269     1. User selects the variogram direction and directional tolerances to apply (horizontal,
270         vertical, and related bandwidth). A maximum separation distance can also be
271         applied to speed up the pair formation process.
272     2. User also selects the required number of variogram points – this controls the
273         number of clusters that will be used by the k-means algorithm.
274     3. Composites database is scanned and composites pairs are formed and stored to a
275         file (Table 1). The 3D separation distance, squared difference of composites grades
276         (semi), and squared difference divided by pair mean (pairwise) are also stored.
277     4. An SPSS syntax file is generated referencing the pairs file.
278     5. SPSS is called using the syntax file and k-means clustering is performed. The Lava
279         script waits for SPSS to complete the clustering process before it continues. The
280         steps performed by SPSS are the following:
281         i.      Opens the pairs file.
282         ii.     Executes k-means clustering (QUICK CLUSTER) based on pair separation
283                 distance. Two new columns are added to the pairs data table – the resulting
284                 cluster number for each pair (QCL_1), and the distance from the cluster centre
285                 (QCL_2).
286         iii.    Aggregates the resulting table based on cluster numbers (QCL_1).
287         iv.     Calculates average separation distance, maximum distance from cluster
288                 centre (QCL_2_max), sum of squared differences (semi), sum of pairwise
289                 squared differences and number of pairs (N_BREAK) for each cluster.
290         v.      Sorts the aggregated table by average separation distance in ascending order
291                 (Table 2).
292         vi.     Saves the aggregated table to a text file.
293     6. The saved table from SPSS is read by the Lava script and is converted to a Vulcan
294         compatible variogram file that can be opened and displayed in Envisage.

295     Table 1: Part of a pairs data table after clustering with k-means in SPSS (from case study 1,
296     azimuth 90$^o$, plunge 20$^o$).

| pair | sample1 | sample2 | distance | hordev | verdev | semi | pairwise |
|------|---------|---------|----------|--------|--------|------|----------|
| P0 | H-314-16.622 | H-316-11.79 | 6.963741379 | 0.095184207 | 0.358188481 | 1161.650889 | 0.471699959 |
| P1 | H-314-16.622 | H-316-10.8 | 7.551314654 | 0.725616456 | 0.129110504 | 321960.9171 | 3.218301598 |
| P2 | H-315-5.539 | H-314-2.778 | 3.194344221 | 0.441519213 | 0.079610374 | 185731.6932 | 2.380098251 |
| P3 | H-315-5.539 | H-314-1.789 | 3.982700968 | 0.153695571 | 0.239766397 | 3180251.622 | 3.483169808 |
| P4 | H-315-5.539 | H-316-1.89 | 4.498334581 | 0.486475155 | 0.003702789 | 8043661.738 | 3.662662936 |
| P5 | H-315-5.539 | H-314-0.8 | 4.846072843 | 0.761050168 | 0.372952977 | 34481440.97 | 3.831501105 |
| P6 | H-315-5.539 | H-316-0.9 | 5.056332663 | 0.347141651 | 0.247642419 | 295832122.5 | 3.941245206 |
| P7 | SD-08-43 | H-15-52.247 | 6.678770995 | 0.301930063 | 1.11746493 | 9846435.237 | 3.643351118 |
| P8 | SD-08-43 | H-15-53.226 | 6.050389822 | 0.223973326 | 0.492179109 | 10400625 | 3.652345679 |
| P9 | SD-08-43 | H-15-54.205 | 5.523861783 | 0.781942854 | 0.116150433 | 21137751.05 | 3.751232685 |
| P10 | M-21-19.7 | H-112-51.75 | 45.89995157 | 2.255489159 | 0.611567557 | 1745181.029 | 3.191502971 |
| P11 | M-21-19.7 | H-112-52.757 | 46.57318349 | 1.510416662 | 0.821022109 | 1745181.029 | 3.191502971 |
| P12 | M-21-19.7 | H-112-53.764 | 47.25828136 | 0.761767907 | 1.026430026 | 1745181.029 | 3.191502971 |

| P13 | M-21-19.7 | H-112-54.771 | 47.95348424 | 0.011819937 | 1.227332402 | 1745181.029 | 3.191502971 |
| P14 | M-21-19.7 | H-112-55.779 | 48.66086482 | 0.7403636 | 1.425085903 | 4041768.472 | 3.438658745 |
| P15 | M-21-19.7 | H-112-56.786 | 49.37828845 | 1.494646484 | 1.618298387 | 9740971.829 | 3.624044192 |
| P16 | M-21-19.7 | H-112-57.793 | 50.10516554 | 2.252006088 | 1.80931579 | 9740971.829 | 3.624044192 |
| P17 | SD-08-42 | H-15-54.205 | 5.46253659 | 0.493360014 | 0.837717765 | 20908497 | 3.671776097 |
| P18 | SD-08-44 | H-15-52.247 | 6.677113448 | 0.05490167 | 0.150186911 | 9609398.609 | 3.474857844 |
| P19 | SD-08-44 | H-15-53.226 | 6.15910586 | 0.506035101 | 0.467754473 | 10156969 | 3.487799746 |
| P20 | SD-08-44 | H-15-51.268 | 7.291150869 | 0.580861371 | 0.782508272 | 10384403.13 | 3.492891538 |

297

298  Table 2: Aggregates table sorted by cluster average separation distance in SPSS (from case
299  study 1, azimuth 90$^o$, plunge 20$^o$).

| QCL_1 | distance_mean | QCL_2_max | semi_sum | pairwise_sum | N_BREAK | semivariogram | pairwise |
|---|---|---|---|---|---|---|---|
| 3 | 8.966 | 8.368 | 4.93E+10 | 1755.058 | 2382 | 2.07E+07 | 0.737 |
| 2 | 19.581 | 5.389 | 6.50E+10 | 2308.253 | 2595 | 2.50E+07 | 0.890 |
| 4 | 30.360 | 8.892 | 4.04E+10 | 1515.350 | 1446 | 2.80E+07 | 1.048 |
| 1 | 48.192 | 8.909 | 6.94E+10 | 2572.557 | 2420 | 2.87E+07 | 1.063 |
| 6 | 63.216 | 13.971 | 5.64E+10 | 1825.052 | 1865 | 3.03E+07 | 0.979 |
| 5 | 91.347 | 14.005 | 1.43E+09 | 105.705 | 222 | 6.46E+06 | 0.476 |

# Case Studies

301  Two case studies using data from real deposits are presented in this paper, selected from a
302  number of examples used to test the VLV approach. Due to the sensitivity of the data, the
303  discussion on its origin and characteristics is kept at a minimum.

304      Experimental variograms were calculated using standard fixed lag variography in
305  Maptek Vulcan 3D, and using variable lag variography with the developed script and IBM
306  SPSS Statistics. The comparison was made on two variography modes, semivariograms and
307  pairwise relative variograms as these are the only modes currently supported by the script.
308  An effort was made to calculate experimental variogram points using both approaches up to
309  the same maximum distance and for the same number of points to make the comparison
310  easier and more objective.

311      A small number of directions were selected in each case to calculate variograms.
312  There are minor differences in the application of azimuth/plunge tolerances and bandwidths
313  in the two approaches compared, as the script does not necessarily replicate the way Vulcan
314  3D forms and selects pairs from particular directions. This has some effect to the difference
315  in the number of pairs reported by the two approaches.

316      The variogram cloud for each direction considered was constructed to gain some
317  understanding for the produced experimental variogram points. The variogram cloud is
318  commonly used as a diagnostic tool in geostatistics and can help detect the presence of
319  outlier points affecting the produced values in the calculation of experimental variograms
320  (Chauvet, 1982, Isaaks et al. 1989, Cressie, 1991). It is essentially a scatter plot of sample pair
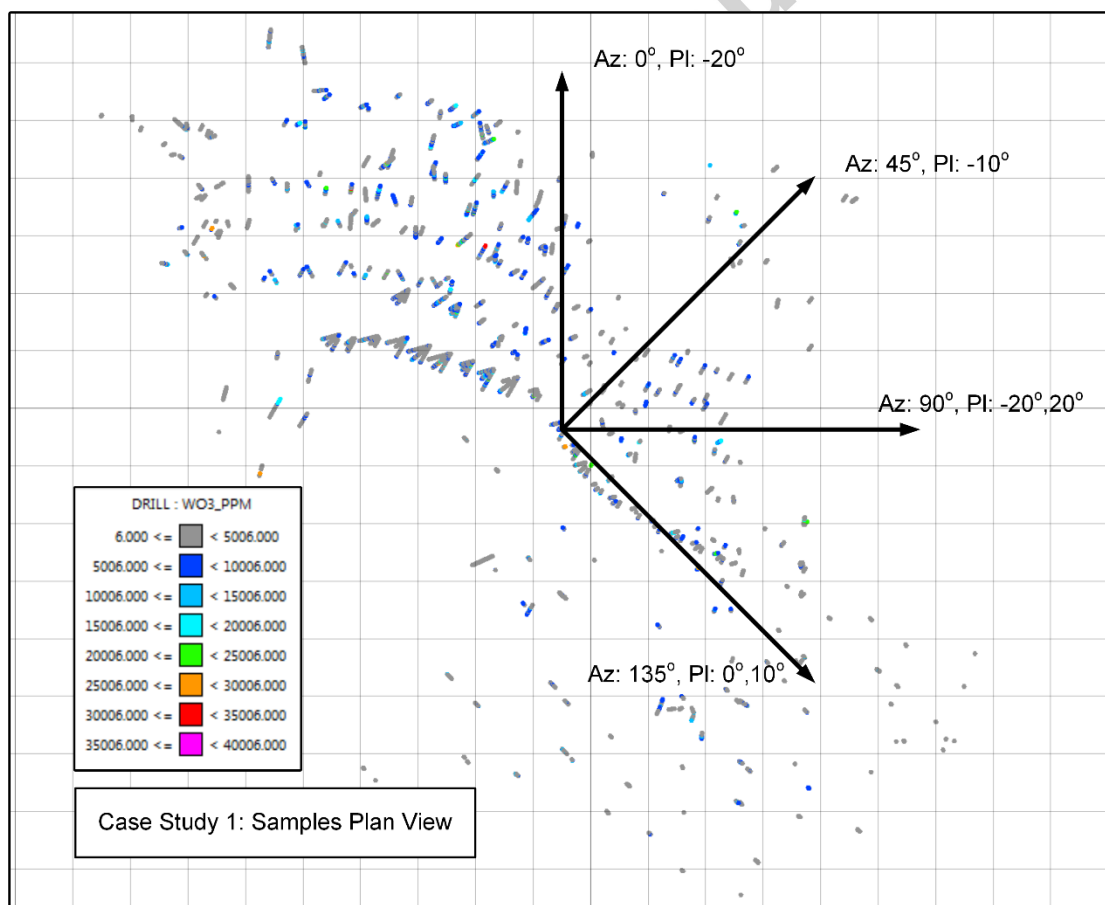321  squared differences against their separation distance (Figures 8 and 15). Other forms of

322 variogram clouds have been proposed and used, such as the square-root differences cloud
323 (Cressie, 1991). The variogram cloud can be used to detect (Plonner, 1999):

324  • Global outliers which are clearly away from the main group of the data.
325  • Local outliers which are more difficult to trace but differ from their neighbouring
326   values and result in high squared differences for small distances.
327  • Small areas of non-stationarity in cases were a cluster of points presents a larger
328   variability than surrounding points.

## Case Study 1 – Tungsten Deposit

330 Data for the first case study come from a tungsten deposit contained within a number of
331 tabular, bedding-conformable skarn horizons. Data includes both underground and surface
332 drilling. The mountainous terrain of the area and the geometry of the underground workings
333 resulted in a fairly irregular sampling pattern. Drillhole samples from one of the tabular
334 horizons were used to produce equal length (1m) composites of $WO_3$ values. A total of 5,466
335 composites were produced and formed the basis for the first case study. Figure 6 shows the
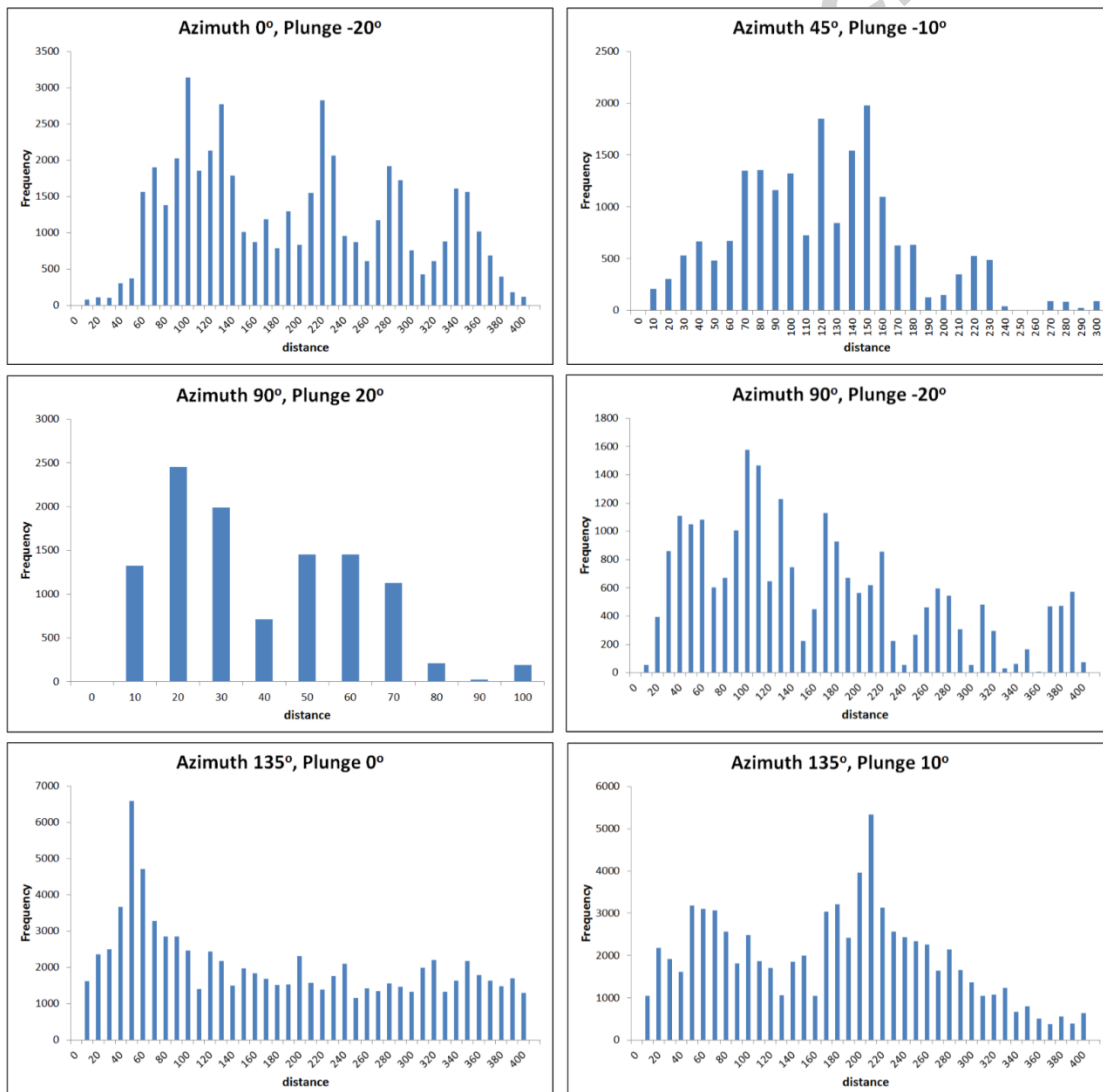336 spatial distribution of these composites in plan view.

337  A number of directions were selected to calculate the experimental variogram in
338 semivariogram and pairwise relative mode. These are shown in Figure 6. At 90° and 135°
339 azimuth, two different plunges were tried, 20° and -20°, and 0° and 10° respectively.



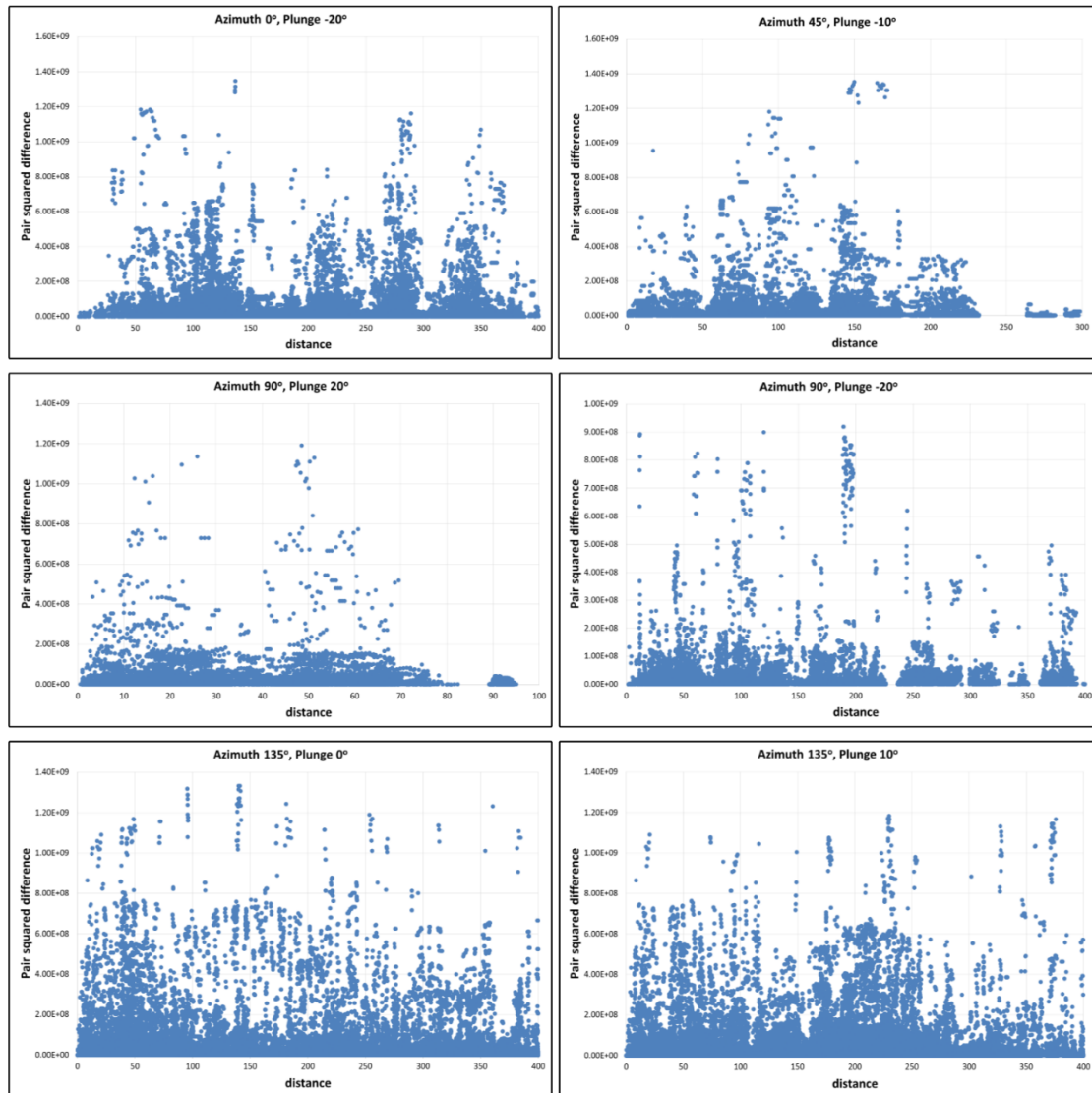Case Study 1: Samples Plan View

340

341 Figure 6: Plan view of underground drillhole sample composites from a particular zone of a
342 tungsten deposit. Irregularity of the pattern is mostly due to drilling following underground
343 openings and drillholes fanning out from almost the same collar location.

344 For each of the directions, a histogram of selected pairs was produced based on the
345 separation distance to help decide the number of clusters for k-means clustering, i.e. the
346 number of variogram points. The same number of points was used (more or less) in standard
347 fixed lag variography to make comparison of the two approaches easier and more
348 conclusive. The histograms were examined visually in order to establish groups of pairs
349 around frequency peaks. In some cases this was possible while in others not quite. For
350 example, in the middle right histogram (Azimuth 90$^o$, Plunge -20$^o$) of Figure 7, the number of
351 peaks is approximately eight. Thus, during k-means clustering, the number of required
352 clusters was set to eight. However, other histograms were not as clear and presented a
353 much more continuous distribution of pairs along the separation distance bins. In those
354 cases, the number of clusters was set according to the number of drillholes along the
355 particular direction up to the maximum distance considered. Clearly, a direct improvement
356 to the VLV approach would be a different clustering algorithm that can set the number of
357 clusters automatically using some criteria (not necessarily just the pair separation distance).



358

359  Figure 7: Histograms of pairs based on separation distance for each of the six directions
360  considered in case study 1.

361  The variogram clouds presented in Figure 8 show the presence of a considerable amount of
362  outliers affecting most distances in each direction considered in this study. The pattern of
363  outliers is mostly uniform in all directions, with the exception of the largest distances of
364  $45°/-10°$ and $90°/20°$ where they are absent. This was reflected in the last couple of points of
365  the corresponding fixed lag variogram graphs (Figure 9) and the last point of the
366  corresponding variable lag variogram graphs (Figure 10).

367



368  Figure 8: Variogram clouds showing the distribution of pair squared difference values along
369  distance for the six directions considered in case study 1.

370

371  Standard fixed lag variography and VLV were performed in all six directions. The lag setups
372  used by the two approaches for each direction are summarised in Table 3. The number of
373  pairs found for each point is also included in the table. In most cases, VLV achieves a much
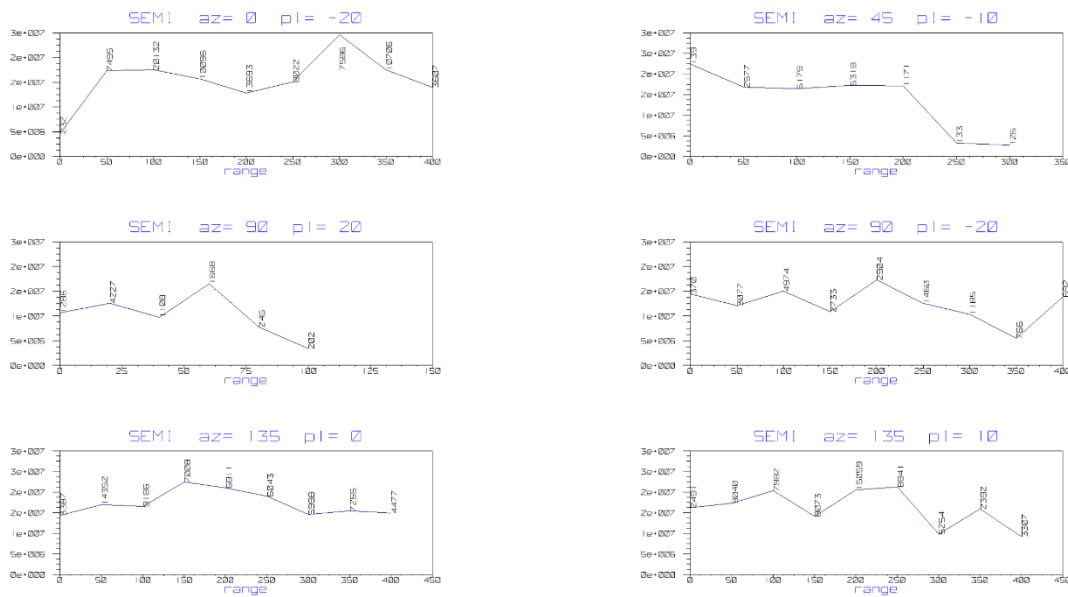
374  more balanced distribution of pairs along the various points, even at higher separation

375  distances.

376  Table 3: Variable lag setup defined by k-means clustering and fixed lag setup defined

377  manually for each of the four directions of case study 1.

| | | Variable Lag | | | | Fixed Lag | | |
|---|---|---|---|---|---|---|---|---|
| | Point | Lag | Lag Tolerance | Number of Pairs | Point | Lag | Lag Tolerance | Number of Pairs |
| Azimuth 0, Plunge -20 | 1 | 57.73 | 56.67 | 5262 | 1 | 50 | 20 | 7495 |
| | 2 | 93.70 | 17.98 | 7702 | 2 | 100 | 20 | 20132 |
| | 3 | 128.34 | 23.98 | 7914 | 3 | 150 | 20 | 10096 |
| | 4 | 176.31 | 23.98 | 4585 | 4 | 200 | 20 | 13693 |
| | 5 | 220.41 | 28.15 | 8274 | 5 | 250 | 20 | 8022 |
| | 6 | 276.74 | 28.13 | 6384 | 6 | 300 | 20 | 7586 |
| | 7 | 330.01 | 26.58 | 4363 | 7 | 350 | 20 | 10706 |
| | 8 | 361.65 | 38.34 | 3027 | 8 | 400 | 20 | 3607 |
| Azimuth 45, Plunge -10 | 1 | 29.29 | 28.29 | 2211 | 1 | 50 | 20 | 2577 |
| | 2 | 71.97 | 21.32 | 4607 | 2 | 100 | 20 | 5179 |
| | 3 | 109.69 | 20.17 | 4651 | 3 | 150 | 20 | 5318 |
| | 4 | 150.18 | 31.45 | 5932 | 4 | 200 | 20 | 1171 |
| | 5 | 213.09 | 31.43 | 1623 | 5 | 250 | 20 | 133 |
| | 6 | 278.82 | 19.90 | 288 | 6 | 300 | 20 | 126 |
| Azimuth 45, Plunge -10 | 1 | 8.97 | 8.37 | 2382 | 1 | 20 | 10 | 4227 |
| | 2 | 19.58 | 5.39 | 2595 | 2 | 40 | 10 | 1108 |
| | 3 | 30.36 | 8.89 | 1446 | 3 | 60 | 10 | 1668 |
| | 4 | 48.19 | 8.91 | 2420 | 4 | 80 | 10 | 246 |
| | 5 | 63.22 | 13.97 | 1865 | 5 | 100 | 10 | 202 |
| | 6 | 91.35 | 14.01 | 222 | 6 | 120 | 10 | 0 |
| Azimuth 90, Plunge -20 | 1 | 39.78 | 37.93 | 4834 | 1 | 50 | 20 | 3077 |
| | 2 | 90.89 | 25.55 | 4802 | 2 | 100 | 20 | 4974 |
| | 3 | 125.06 | 22.41 | 3007 | 3 | 150 | 20 | 2733 |
| | 4 | 169.95 | 22.43 | 3193 | 4 | 200 | 20 | 2904 |
| | 5 | 207.97 | 18.99 | 2319 | 5 | 250 | 20 | 1460 |
| | 6 | 264.20 | 27.31 | 2197 | 6 | 300 | 20 | 1185 |
| | 7 | 310.20 | 31.72 | 958 | 7 | 350 | 20 | 766 |
| | 8 | 373.64 | 31.51 | 1748 | 8 | 400 | 20 | 692 |
| Azimuth 135, Plunge 0 | 1 | 24.26 | 24.10 | 10470 | 1 | 50 | 20 | 14352 |
| | 2 | 57.04 | 23.74 | 17363 | 2 | 100 | 20 | 9186 |
| | 3 | 104.54 | 26.54 | 11232 | 3 | 150 | 20 | 7008 |
| | 4 | 157.62 | 28.10 | 9312 | 4 | 200 | 20 | 6811 |
| | 5 | 213.85 | 28.10 | 9926 | 5 | 250 | 20 | 6043 |
| | 6 | 268.69 | 27.42 | 7321 | 6 | 300 | 20 | 6998 |
| | 7 | 319.53 | 25.61 | 9051 | 7 | 350 | 20 | 7255 |
| | 8 | 370.75 | 29.22 | 9025 | 8 | 400 | 20 | 4477 |
| Azimuth 135, Plunge 10 | 1 | 22.48 | 22.33 | 7161 | 1 | 50 | 20 | 9040 |
| | 2 | 60.46 | 22.26 | 11820 | 2 | 100 | 20 | 7982 |
| | 3 | 105.00 | 28.42 | 9247 | 3 | 150 | 20 | 8073 |
| | 4 | 161.86 | 28.42 | 11463 | 4 | 200 | 20 | 15059 |
| | 5 | 205.83 | 23.51 | 16326 | 5 | 250 | 20 | 8841 |
| | 6 | 252.85 | 26.04 | 10796 | 6 | 300 | 20 | 5254 |
| | 7 | 304.95 | 31.46 | 7146 | 7 | 350 | 20 | 2392 |
| | 8 | 367.95 | 32.05 | 3388 | 8 | 400 | 20 | 3307 |

378

379          Figures 9 and 10 present the standard semivariogram in each direction produced by
380   fixed lag variography and VLV respectively. It must be noted that lag tolerances in the case
381   of fixed lag variography were set after some testing while VLV lag tolerances were set
382   automatically by the clustering process. Figures 11 and 12 present the pairwise relative
383   variogram in the same manner. In both variogram modes and in most cases, the points
384   produced by VLV define a smoother, easier to interpret, graph.

385



386

387   Figure 9: Standard semivariograms produced using fixed lag variography for case study 1.
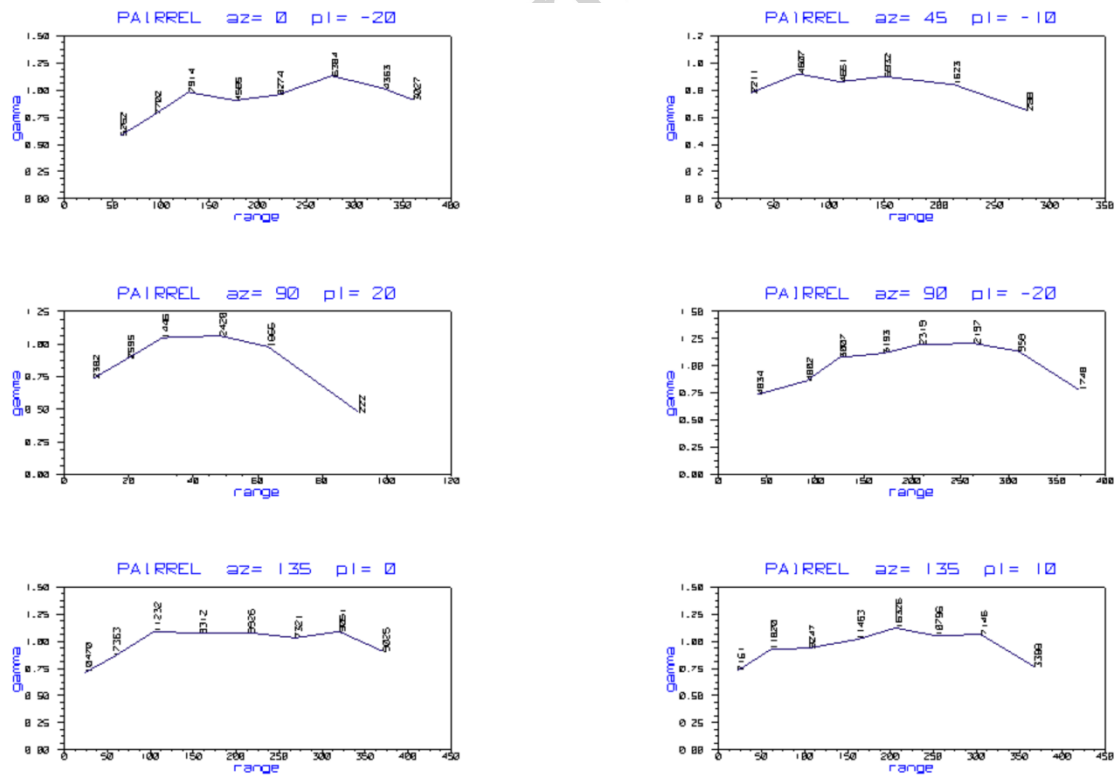
388



389

390    Figure 10: Standard semivariograms produced using variable lag variography for case study
391    1.



392

393    Figure 11: Pairwise relative variograms produced using fixed lag variography for case study
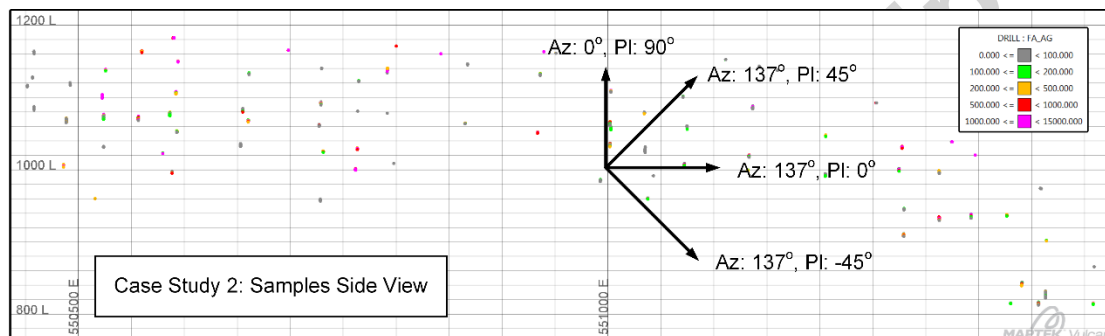394    1.



395

396 Figure 12: Pairwise relative variograms produced using variable lag variography tools for
397 case study 1.

398

399        Running the Perl script to perform VLV required considerably more time than
400 running normal variography tools in Vulcan (a few minutes compared to a few seconds). This
401 is due to the fact that Perl is an interpreted language, and pairs were written to and read
402 from an ASCII file to save memory.

## Case Study 2 – Silver Vein Deposit

404 Data for the second case study include 573 composites of approximately 1m length from a
405 near vertical silver vein. Two separate drilling campaigns (original exploration plus infill
406 drilling) resulted in some irregularity of the sampling pattern – considerably less though
407 compared to the first case study. Figure 12 shows a side view of the data and the selected
408 directions for experimental variogram calculation.

409


410 Figure 13: Side view of samples from a vertical silver vein showing a fair amount of
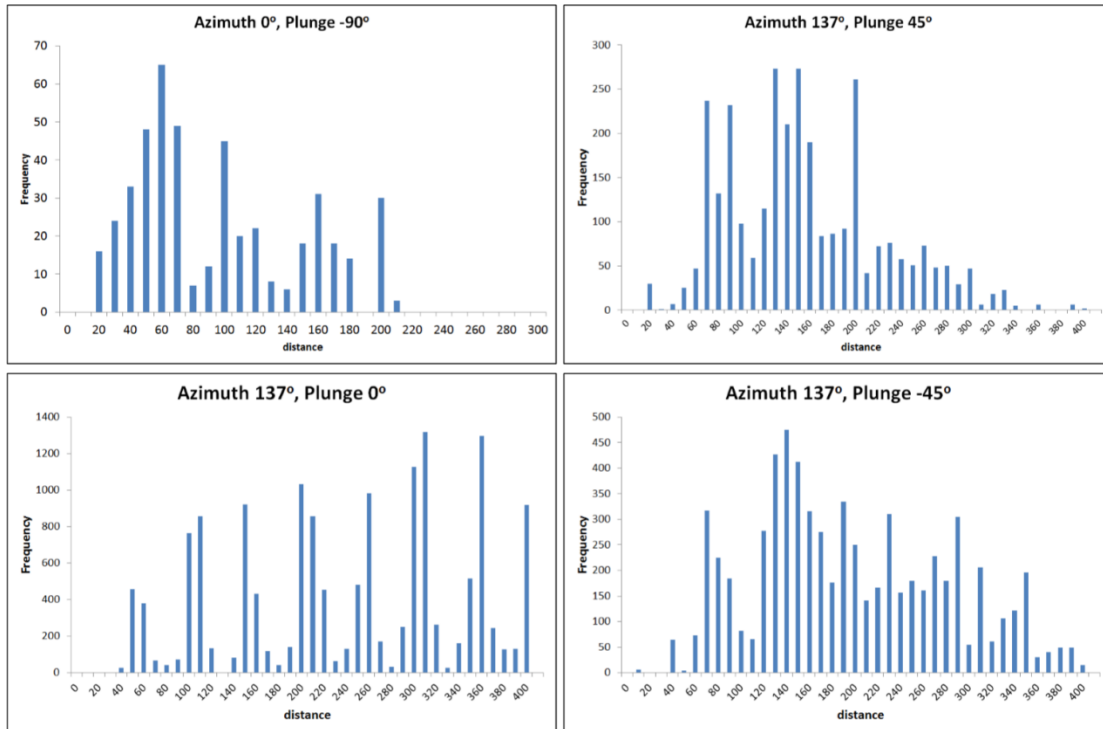411 irregularity in spacing, partly due to infill drilling.

412

413 Histograms were constructed with the pairs according to separation distance as in case
414 study 1. These reflected the much more regular sampling pattern (Figure 14). Particularly the
415 one along the strike direction (azimuth 137$^o$, plunge 0$^o$) presents exactly the original drillhole
416 spacing (50m) through the corresponding peaks. Thus, in the second case study, it was much
417 easier to decide the number of required clusters or variogram points to calculate.

418        The lag setup used in fixed lag variography and the setup configured with VLV are
419 shown in Table 4. The differences in the way directional tolerances are applied in the two
420 approaches were quite evident in the number of pairs reported. VLV still produced a slightly
421 more balanced distribution of pairs on the different variogram points, but not to the extent
422 shown in the first case study as the sampling pattern is much more regular this time and it is
423 much easier for a fixed lag setup to follow it.

424        The variogram clouds for the four directions considered are shown in Figure 15. In
425 this case study, the distribution of outliers is not as uniform across the range of distances
426 (with the exception of direction 137$^o$/0$^o$ which follows the drilling pattern more closely). This
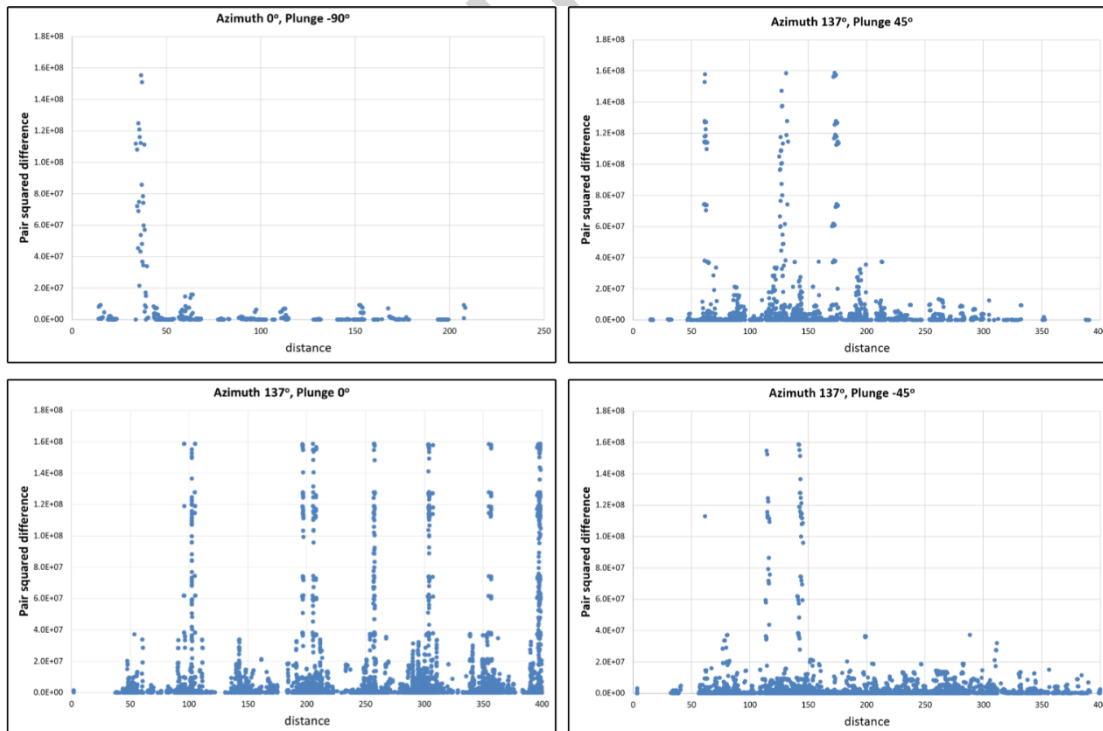
427 means that there could be room for improvement of the semivariogram graphs by trimming
428 of outliers at appropriate levels.



429

430 Figure 14: Histograms of pairs based on separation distance for each of the four directions
431 considered in case study 2.

432



433

434 Figure 15: Variogram clouds showing the distribution of pair squared difference values along
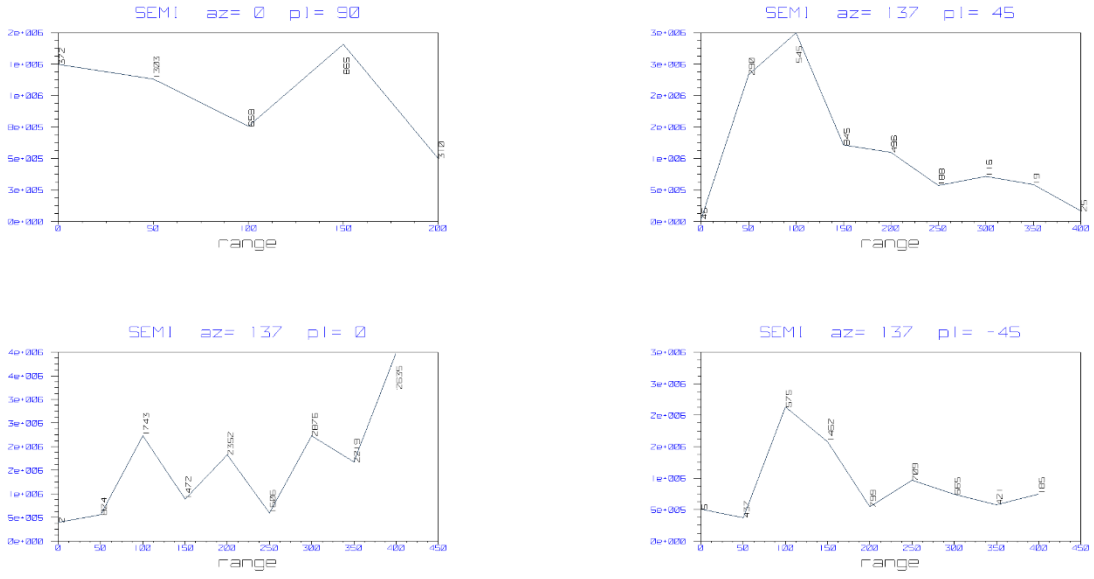435 distance for the four directions considered in case study 2.

436 Table 4: Variable lag setup defined by k-means clustering and fixed lag setup defined
437 manually for each of the four directions of case study 2.

| | | Variable Lag | | | | Fixed Lag | | |
|---|---|---|---|---|---|---|---|---|
| | Point | Lag | Lag Tolerance | Number of Pairs | Point | Lag | Lag Tolerance | Number of Pairs |
| Azimuth 0, Plunge -90 | 1 | 28.12 | 13.47 | 70 | 1 | 50 | 15 | 1303 |
| | 2 | 56.72 | 23.33 | 156 | 2 | 100 | 15 | 658 |
| | 3 | 103.61 | 28.14 | 105 | 3 | 150 | 15 | 865 |
| | 4 | 169.33 | 30.07 | 104 | 4 | 200 | 15 | 310 |
| Azimuth 137, Plunge 45 | 1 | 36.31 | 21.33 | 82 | 1 | 50 | 20 | 290 |
| | 2 | 78.45 | 28.70 | 772 | 2 | 100 | 20 | 545 |
| | 3 | 136.14 | 28.75 | 1091 | 3 | 150 | 20 | 845 |
| | 4 | 190.22 | 28.42 | 617 | 4 | 200 | 20 | 496 |
| | 5 | 247.18 | 30.99 | 359 | 5 | 250 | 20 | 188 |
| | 6 | 309.48 | 81.06 | 143 | 6 | 300 | 20 | 116 |
| | | | | | 7 | 350 | 20 | 19 |
| | | | | | 8 | 400 | 20 | 25 |
| Azimuth 137, Plunge 0 | 1 | 51.22 | 49.31 | 941 | 1 | 50 | 20 | 824 |
| | 2 | 99.86 | 23.91 | 1862 | 2 | 100 | 20 | 1743 |
| | 3 | 149.26 | 26.03 | 1598 | 3 | 150 | 20 | 1472 |
| | 4 | 201.66 | 24.33 | 2536 | 4 | 200 | 20 | 2352 |
| | 5 | 251.72 | 23.26 | 1798 | 5 | 250 | 20 | 1606 |
| | 6 | 300.28 | 24.30 | 2991 | 6 | 300 | 20 | 2876 |
| | 7 | 352.94 | 24.47 | 2294 | 7 | 350 | 20 | 2219 |
| | 8 | 394.03 | 20.42 | 1111 | 8 | 400 | 20 | 2635 |
| Azimuth 137, Plunge -45 | 1 | 70.52 | 66.98 | 968 | 1 | 50 | 20 | 437 |
| | 2 | 134.33 | 31.82 | 1923 | 2 | 100 | 20 | 575 |
| | 3 | 180.65 | 25.21 | 1132 | 3 | 150 | 20 | 1462 |
| | 4 | 231.15 | 26.24 | 1038 | 4 | 200 | 20 | 799 |
| | 5 | 283.99 | 31.32 | 1029 | 5 | 250 | 20 | 709 |
| | 6 | 346.71 | 53.28 | 630 | 6 | 300 | 20 | 665 |
| | | | | | 7 | 350 | 20 | 421 |
| | | | | | 8 | 400 | 20 | 185 |

438

439       Figures 16 and 17 show the standard semivariogram produced by both approaches.
440 The benefit of using VLV over fixed lag variography is evident as the produced points define
441 a very clear structure. In pairwise relative mode (Figures 18 and 19) the improvement is
442 smaller but still significant. The time required to run VLV was considerably less compared to
443 the first case study as the number of composites and possible pairs was much smaller. The
444 actual clustering process in IBM SPSS Statistics required a couple of seconds in both case
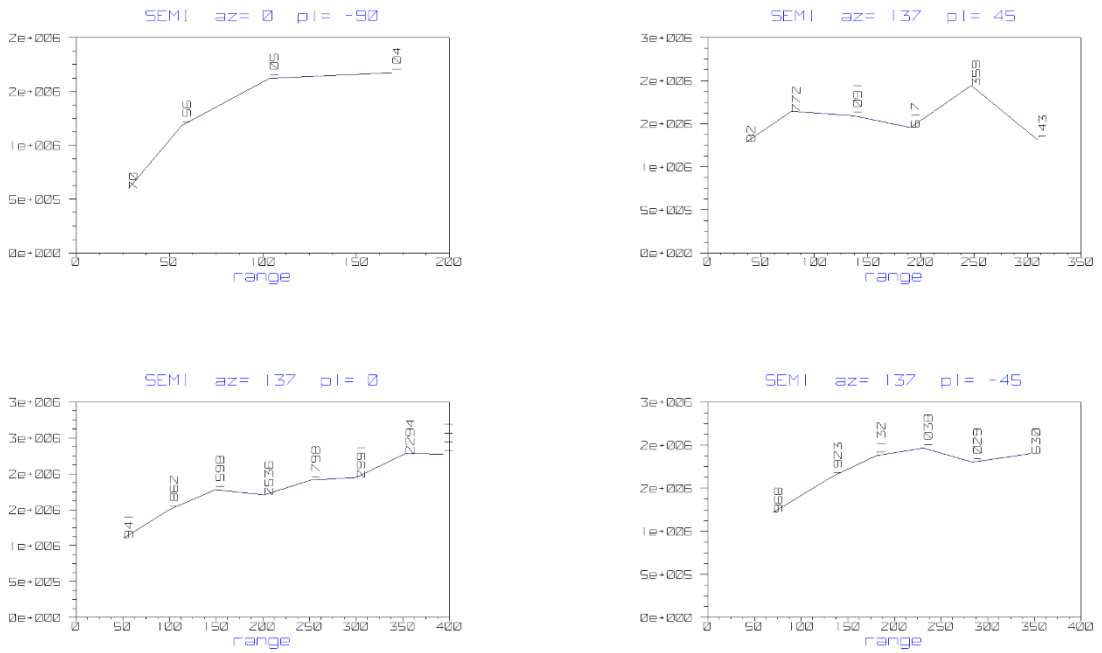445 studies for each of the directions.

446

447

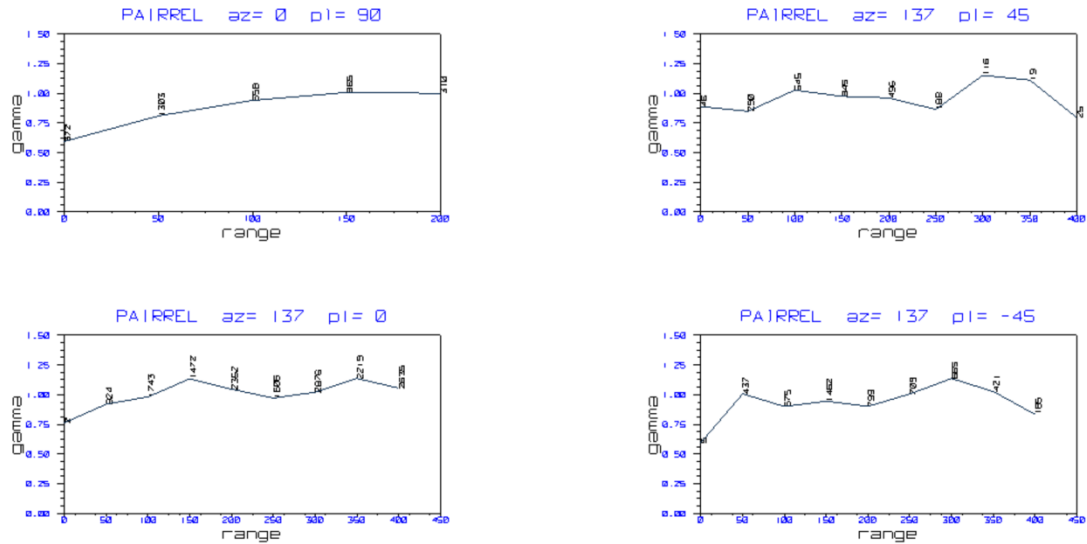448     Figure 16: Standard semivariograms produced using fixed lag variography for case study 2.

449



450

451     Figure 17: Standard semivariograms produced using variable lag variography for case study
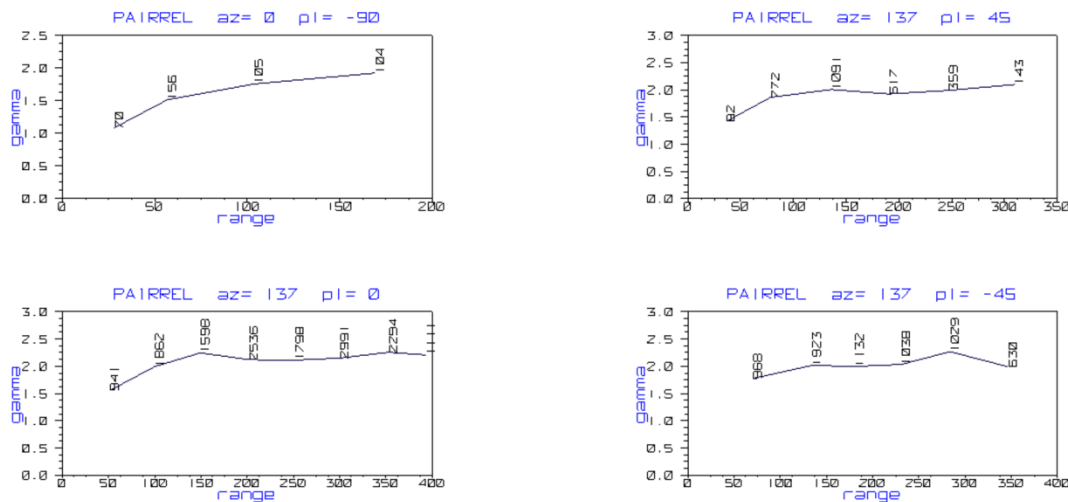452     2.

453

454

Figure 18: Pairwise relative variograms produced using fixed lag variography for case study 2.



Figure 19: Pairwise relative variograms produced using variable lag variography tools for case study 2.

## Conclusions

This paper presented an alternative method for calculating experimental variograms based on variable lags defined through a clustering process. The development of this method was motivated by the excessive time and effort required in setting up lags and tolerances for variography in cases of irregular sampling patterns. K-means clustering was considered as the algorithm for clustering sample pairs based on separation distance. A script was developed that runs through a mine planning package, which creates the pairs, runs the clustering process through a statistical software package and produces the experimental variogram in two variography modes (semivariogram and pairwise relative). The method was tested on data from a number of real deposits, two of which are presented as case

471 studies in this paper. Normal variography was also performed using standard geostatistical
472 tools in order to evaluate the benefits of the variable lag variography (VLV) concept. An
473 effort was made to keep all variography parameters (other than lag and lag tolerance) the
474 same to make the comparison easier and more effective.

475      The results from calculating experimental variograms using both approaches have
476 shown that VLV can relieve the practitioner from the trouble of finding an appropriate lag
477 setup and at the same time produce experimental variograms which are smoother and
478 easier to interpret in cases of irregular sampling patterns. More testing is required with
479 other datasets to have a better understanding of the effects of using VLV.

480      Some input is still required to VLV as the k-means algorithm does not define the
481 number of clusters automatically. It is one of the aims for future work to develop a method
482 to find the optimum number of clusters (or variogram points) based on a given set of pairs.
483 Currently, VLV is using separation distance as the sole criterion for clustering. Other criteria
484 are considered, such as the squared difference of the pair sample grades. The effect of such
485 criteria needs to be evaluated.

486      The script that was developed to perform VLV will also be rewritten to speed up the
487 pair formation process. An implementation of the k-means algorithm will also be included in
488 the scrip so that an external statistical package is no longer required, if such development
489 does not have a negative effect on the speed of clustering.

490 ## References

491 Bleines, C. J., Deraisme, F., Geffory, N., Jeannee, S., Perseval, F., Rambert, D., Renard, O.,
492 Torres, and Touffait, Y. (2013) Isatis Beginner's Guide, Geovariances & Ecole Des Mines De
493 Paris.

494 Chauvet, P. (1982). The Variogram Cloud. In Proceedings of the 17th APCOM Symposium.
495 Colorado School of Mines, Golden, CO, pages 757-764.

496 Cressie, N. A. C. (1991). Statistics for Spatial Data. New York: Wiley & Sons.

497 Deutsch, C.V., and Journel, A.G. (1992) GSLIB - Geostatistical Software Library and User's
498 Guide, Oxford University Press.

499 Englund, E., and Sparks, A. (1991). GEO-EAS 1.2.1 User's Guide, US Environmental Protection
500 Agency & Computer Sciences Corporation.

501 Forgy, E.W. (1965). Cluster analysis of multivariate data: efficiency versus interpretability of
502 classifications. Biometrics 21: 768–769.

503 Hartigan, J. A. (1975). Clustering algorithms. New York: John Wiley and Sons.

504 Hartigan, J. A.; Wong, M. A. (1979). Algorithm AS 136: A K-Means Clustering Algorithm.
505 Journal of the Royal Statistical Society, Series C 28 (1): 100–108.

506 IBM SPSS Statistics 20 Algorithms (2011), IBM Corporation.

507    IBM SPSS Statistics 20 Base (2011), IBM Corporation.

508    Isaaks, E. H. and Srivastava, R. M. (1989). An Introduction to Applied Geostatistics. New York:
509    Oxford University Press.

510    Lloyd, S. P. (1982). Least square quantization in PCM. Bell Telephone Laboratories Paper.
511    IEEE Transactions on Information Theory 28 (2): 129–137.

512    McQueen, J. B. (1967). Some Methods for classification and Analysis of Multivariate
513    Observations. Proceedings of 5th Berkeley Symposium on Mathematical Statistics and
514    Probability 1. University of California Press. pp. 281–297.

515    Pannatier, Y. (1996). VARIOWIN: Software for Spatial Data Analysis in 2D. Springer-Verlag,
516    New York.

517    Ploner, A. (1999). The Use of the Variogram Cloud in Geostatistical Modelling,
518    Environmetrics, 10, 413-437: John Wiley & Sons.

519    Remy, N., Boucher, A., and Wu, J. (2011). Applied Geostatistics with SGeMS – A User's Guide.
520    Cambridge.