

調査の進行に伴う分散の増加

正路 徹也

Increase of Variance with Progressing Survey

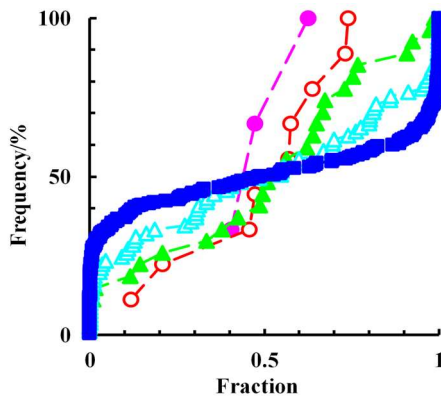
Tetsuya Shoji*

* 東京大学 The University of Tokyo, Tokyo 113-8639, Japan.
E-mail: t-t_shoji@jcom.home.ne.jp

キーワード：地球科学，調査，情報，評価，分散

Key words: Geoscience, Survey, Information, Evaluation, Variance

地球科学的調査の進行に伴い対象地域を構成する地質体の分布およびそれらの接触状態に関する地球情報が増加する。したがって、対象地域を構成する各セルにおける各地質体の存在度の推定値の精度が向上して、エントロピーあるいは分散の値が小さくなる。この調査の進行に伴うエントロピーあるいは分散が減少する様子を、正路(2018)は次の単純な数値モデルで示した。調査は直線上で行われ、段階1では全範囲の midpoint の地質体を同定する。段階2では、全体を3つのセルに分け、各セルの midpoint の地質体を同定する(このうち中間のセルの midpoint は、段階1で調査済み)。段階3では各セルをさらに3つのセルに分け、各セルの midpoint の地質体を同定する。以下、この操作を繰り返し、段階6でセルの大きさが目的の分解能に達し調査は終了する。条件を変えたいくつかのシミュレーションのうち、地質体の種類数が2(他は4)で、存在割合が準二値(他は完全二値)で与えられ、存在割合の比が50:50(他は20:80, 10:90, 5:95, 2:98)の場合について、調査段階をパラメーターとして存在割合の累積頻度を第1図に示す(5本の点列のうち、構成点数が少ない方から段階2, 3, 4, 5, 6)。一見して明らかなように、調査の進行とともに分散が増加している。実際、分散の値



第1図. 地質体の種類数が2で、存在割合が準二値で与えられ、存在割合の比が50:50の場合における存在割合の累積頻度。各点列は調査段階に対応し、構成点数が少ない方から段階2, 3, 4, 5, 6。分散の値は、この順序で0.013, 0.047, 0.096, 0.149, 0.196。

は、この順序で0.013, 0.047, 0.096, 0.149, 0.196と増加している。これは正路(2018)の結論と逆である。その理由は、以下に詳述するように分散の定義の違いによる。

正路(2018)では、1つのセル*i*(総数*N*)に注目し、そこにおける地質体*j*(種類数*m*)の存在割合が p_{ij} の場合、当該セルの分散 $V_{i\text{Indv}}$ は($m-1$)次元空間の正単体 $A_1A_2\cdots A_m$ の頂点 A_j に重み p_{ij} があり p_{ij} ($j=1, \dots, m$)で与えられる単体内の1点 P_i の周りの2次モーメントとして次式で定義する。

$$V_{i\text{Indv}} = \sum_{j=1}^m p_{ij} \left\{ \sum_{j_1=1}^m p_{ij_1}^2 + \sum_{j_1=1}^{m-1} p_{ij_1} \sum_{j_2=j_1+1}^m p_{ij_2} - p_{ji} \sum_{j_1=1}^m p_{ij_1} \right\} \\ = \sum_{j_1=1}^{m-1} p_{ij_1} \sum_{j_2=j_1+1}^m p_{ij_2} \quad (1)$$

ここで、添え字の **Indv**=**Individual** は個別のセルに注目していることを意味し、{}内は点 P_i と頂点 A_j の距離の二乗である。これより全体の分散 V_{Indv} は次式で与えられる。

$$V_{\text{Indv}} = \sum_{i=1}^N V_{i\text{Indv}} / N$$

なお、式(1)で与えられる分散は、セル*i*のデータのみを使って計算されているので、以後個別分散(**individual variance**)と呼ぶ。

二元系($m=2$)の場合、地質体1と2の存在割合をそれぞれ p_{1i} と p_{2i} ($p_{1i}+p_{2i}=1$)とすると、点 P_i と頂点 A_1 との距離は p_{1i} 、頂点 A_2 との距離は p_{2i} であるから、式(1)は次式(2)のように表される。

$$V_{i\text{Indv}} = p_{1i}p_{2i}^2 + p_{2i}p_{1i}^2 = p_{1i}p_{2i}(p_{2i} + p_{1i}) \\ = p_{1i}p_{2i} = p_{1i}(1 - p_{1i}) \quad (2)$$

第1図に示されている累積頻度の分散 V_{Bulk} (添え字 **Bulk**の意味は後述)は、セル*i*の存在割合を p_i (二元系なので地質体の種類を示す添え字は省略)として、次式(3)で与えられる。

$$V_{\text{Bulk}} = \sum_{i=1}^N (p_i - \bar{p})^2 / N \\ = \sum_{i=1}^N \left(p_i - \sum_{i=1}^N p_i / N \right)^2 / N \quad (3)$$

ここで、 \bar{p} は全データ(**bulk data**)の平均である。式(2)に相当するセル*i*のバルク分散 $V_{i\text{Bulk}}$ は次式(4)で与えられる。

$$V_{i\text{Bulk}} = (p_i - \bar{p})^2 \quad (4)$$

式(4)には、 \bar{p} というセル*i*以外のセルの値を必要とする全体の平均が入っている。すなわち、1つのセルのバルク分散を求めるとき、個別分散とは違い、当該セルのみの値では不十分で、他のセルの値も必要とする。これがバルク(bulk)の意味である。

第2図の左右にそれぞれ調査の進行に伴う個別分散とバルク分散の変化を示す。両図を比較すると、1) 個別分散は段階0に値があるが、バルク分散にはない、2) バルク分散は個別分散の上下を反転したパターンを示す。段階0における個別分散が定義できるのは、対象地域全体の地質体の存在割合が最終段階における存在割合の平均に等しいと推定し、その値に相当する重みが地質体を表す両端に掛かっているとするためである。バルク分散は個別分散の上下を反転したパターンについては後述する。

第1図の累積頻度図を微分すると、微分型頻度図が得られる。上記の分散の説明はこの微分型頻度図に基づいている。三元系では、累積する方向によって得られる累積頻度図が変わる。しかし、微分型頻度図を考えれば、分散の計算は可能である。例えば、*m*元系の場合、セル*i*の地質体の存在割合を(*m*-1)次元の正単体内の1点 P_i として表し、その点と平均の存在割合を表す点 P_M との距離 $\overline{P_i P_M}$ の二乗を l_i^2 とすれば、それが次式のようにバルク分散となる。

$$V_{i\text{Bulk}} = l_i^2$$

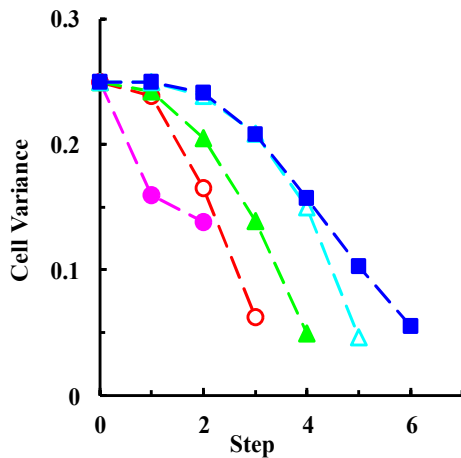
三元系(*m* = 3)で、高さが1の正三角形 $A_1 A_2 A_3$ 内の点 P_i の位置を辺 $A_2 A_3$, $A_1 A_3$, $A_1 A_2$ に下した垂線の長さそれぞれ p_{1i} , p_{2i} , p_{3i} ($p_{1i} + p_{2i} + p_{3i} = 1$)で、点 P_M の位置も同様に \bar{p}_1 , \bar{p}_2 , \bar{p}_3 で表すと、長さ $\overline{P_i P_M}$ の二乗 l_i^2 は次式(5)で与えられ、これがセル*i*のバルク分散となる。

$$V_{i\text{Bulk}} = l_i^2 = \Delta p_1^2 + \Delta p_1 \Delta p_2 + \Delta p_2^2 \quad (5)$$

ここで、 $\Delta p_1 = p_{1i} - \bar{p}_1$, $\Delta p_2 = p_{2i} - \bar{p}_2$ で、

$$\bar{p}_1 = \sum_{i=1}^N p_{1i} / N$$

$$\bar{p}_2 = \sum_{i=1}^N p_{2i} / N$$



である。これより全体のバルク分散 V_{Bulk} は次式で与えられる。

$$V_{\text{Bulk}} = \sum_{i=1}^N V_{i\text{Bulk}} / N$$

前述のよう、個別分散とバルク分散は上下を反転したパターンをなす(第2図)。したがって、両者は負相関を示している。さらに、二元系で両者の和は次式で示すように、 $(\bar{p}_1 - \bar{p}_1^2)$ と一定値である。

$$\begin{aligned} V_{\text{Indv}} + V_{\text{Bulk}} &= \left\{ \bar{p}_1 - \sum_{i=1}^N p_{1i}^2 / N \right\} \\ &\quad + \left\{ \sum_{i=1}^N p_{1i}^2 / N - \bar{p}_1^2 \right\} \\ &= \bar{p}_1 - \bar{p}_1^2 \end{aligned}$$

すなわち、両者は完全負相関($r = -1$)をなす。三元系においても、次式に示すように完全負相関をなす。

$$\begin{aligned} V_{\text{Indv}} + V_{\text{Bulk}} &= (\bar{p}_1 + \bar{p}_2) \\ &\quad - (\bar{p}_1^2 + \bar{p}_1 \bar{p}_2 + \bar{p}_2^2) \end{aligned}$$

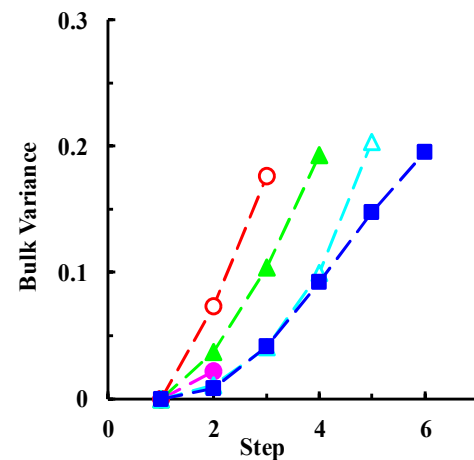
しかし、セル*i*の個別分散とバルク分散は次式のような関係にあり、負相関はなすが、完全負相関ではない。

$$V_{i\text{Indv}} + V_{i\text{Bulk}} = p_{1i} - 2p_{1i}\bar{p}_1 + \bar{p}_1^2$$

以上より次のまとめが得られる。正路(2018)は、調査の進行に伴う分散の減少の様子を簡単な数値モデルで示した。しかし、通常の統計処理で用いられている分散を使うと、地質体の存在割合の分散は調査の進行に伴って増加する。前者の分散は個別セルの値のみで求められるので個別分散、後者の分散を求めるために全データの平均が必要なのでバルク分散と呼ぶと、両者は完全負相関の関係にある。

引用文献

正路徹也(2018):地球科学的調査の進行に伴う地球情報エントロピーと分散の減少. 情報地質, 29(2), 61-75.



第2図. 調査の進行に伴うセル分散(左)とバルク分散(右)の変化. 各点列の右端が最終段階に対応する。