連続値を対象とした位相的階層構造に基づく 空間集積性の検出について

石岡 文生†

(受付 2024 年 6 月 24 日;改訂 10 月 29 日;採択 10 月 31 日)

要 旨

「興味のある対象が,ある特定の地域に集中しているか(空間的な集積性は存在するか)」を知 ることは、空間データ解析における関心事の一つである.近年、空間データの各領域を特定の ルールに基づいて走査(スキャン)し、尤度に基づいて空間集積性の有無を評価する空間スキャ ン検定が、様々な分野において広く利用されている.しかし、空間スキャン検定に関連する研 究成果の多くは、主にその対象をカウントデータ(離散値)としており、また、その集積領域群 の形状にも制約があった.そこで本研究では、「離散値をとらない空間データに対し、任意の 形状をした集積領域群を評価できるか」という問いに対し、エシェロン解析法をベースとした アプローチによって解決を試みる方法を提案する.エシェロン解析法とは、空間データの各領 域が持つ1変量値と領域間の近傍情報に基づいて、データを同位相の領域群(エシェロン)に分 類し、それらを階層構造のグラフで表現する手法である.提案法により、多くの場合で連続値 として得られる「ある予測モデルに基づいた推定値」などのデータに対して、柔軟な空間集積性 の議論が可能となる.本稿では、数値実験を通じて提案法の有効性を検証するとともに、従来 法との違いについて考察する.また、クリギング予測値やベイズ推定値といった連続値を取る データへの応用例も紹介する.

キーワード:空間集積性,空間スキャン検定,エシェロン解析法,エシェロンスキャン法.

1. はじめに

「興味のある対象が,ある特定の地域に集中しているか(空間的な集積性は存在するか)」を知 ることは、データ解析における関心事の一つである。例えば、「新型肺炎などの感染症がどの 地域で集中して発生しているのか?」や、「ある町の生活道路において、交通事故がどの場所で 頻発しているのか?」などのように、「どこかに集積性は存在するのか?それとも全体的にばら ついているのか?」「集積性が存在しているとしたら、その範囲はどこまでなのか?」を統計的 根拠(データ)に基づいて決定できれば、「発生源の特定による原因の解明」や、「将来の環境や 健康に対する影響の早期発見」に寄与することが期待される。

そんな中,空間データの各領域を特定のルールに基づいて走査(スキャン)し,尤度に基づい て空間集積性の有無を評価する空間スキャン検定 (Kulldorff, 1997)が,疫学の分野などで広く 利用されている.また,この空間スキャン検定を行うためのフリーソフトウェア SaTScan™ と 相まって,関連する研究成果が多数提供されている.空間スキャン検定は,1. 尤度算出に利用する統計モデル,および2. スキャン方式,の2つの要素から成り立ち,それぞれの要素を組み合わせることで実行される.ところが,実際の空間スキャン検定の適用において,1の側面では「ある疾病に罹患した患者数」や「交通事故の発生件数」などのカウントデータ(離散値)を対象にした Poisson モデルを用いたもの,2の側面では領域を同心円状にスキャンする circular scan 法を用いたもの,がその大半を占める.そのため,「離散値をとらないデータ」を対象とした「任意形状をした空間集積性」を評価するための空間スキャン検定は,十分に確立されていないのが現状である.

離散値をとらないデータに関する空間集積性の先行研究では、連続値を対象とした Normal モデル (Kulldorff et al., 2009; Huang et al., 2009)や、生存時間を対象にした Exponential モデ ル (Huang et al., 2007)およびそれを拡張した Weibull モデル (Bhatt and Tiwari, 2014)などが 提案されている. これらのモデルの一部は、ソフトウェア SaTScan™ に実装されているが、こ のソフトは同心円状にスキャンを行う circular scan 法 (Kulldorff, 1997)を採用しているため、 円の形状をした集積領域群しか検出できず、例えば川や道路に沿うような非円状の集積性は評 価できない. SaTScan™ には楕円状にスキャンするオプションも存在するが、これはスキャン 形状の課題を本質的に解決するものにはなっていない. 一方で、検出される集積形状の制約を 緩和するスキャン方式が複数提案されており (Duczmal and Assunção, 2004; Patil and Taillie, 2004; Tango and Takahashi, 2005; Costa et al., 2012)、そのいくつかは統計ソフト R のパッケー ジで公開されている. しかしながら、それらの多くは Poisson モデルや Bernoulli モデルなどの 離散値を想定したモデルでデザインされている.

このような状況の中で,筆者らはエシェロンスキャン法(栗原,2003;石岡・栗原,2012)と 呼ばれるスキャン方式を提案している.筆者らは、この方式によって「集積形状の制約」や「計 算コスト」が改善されることを実証してきた(石岡・栗原,2012; Ishioka et al., 2019; 竹村 他, 2021).一方で,これら先行研究では主に「領域のスキャンの効率化」に焦点を当てており,他 のスキャン方式との比較検証や、実データへの適用の際には、離散値を用いてきた.連続値を 対象にエシェロンスキャン法を適用した先行研究としては、掃部 他(2023)があるものの、そ の有効性についての詳細な検証や応用についてはまだ検討の余地が残されている.そこで本稿 では、空間スキャン検定における重要な課題の一つ、すなわち「離散値をとらない空間データ に対し、任意の形状をした集積領域群を評価できるか」について、エシェロンスキャン法のア プローチから解決を試みる.

本稿の構成は以下の通りである。第2節おいて本稿で扱う統計量のモデルについて述べる。 第3節ではエシェロンスキャン法について説明する。第4節では提案法の有効性について数値 実験により検証する。第5節では、提案法の応用例としてクリギング予測値やベイズ推定値へ の適用を試みる。第6節でまとめと今後の課題について触れる。

2. 空間スキャン統計量

2.1 連続値を対象とした Normal モデル

空間スキャン統計量は、空間集積性を評価するのに用いられる統計指標であり、最初は Kulldorff and Nagarwalla (1995)によって提案された. その後、データの性質に応じた様々な 統計量のモデルが提案されているが、本稿では、連続的なデータに適応するために開発された Normal モデル (Kulldorff et al., 2009; Huang et al., 2009)を扱う. これは、対象データが個々 のカウント値として取得されるのではなく、地域ごとに集計された感染率や汚染濃度などのよ うに、個別の観測値を基に推定される集約レベルの観測値に対して、未知の真値の地理的分布 に焦点を当てたモデルとなっている.

いま,解析対象地域全体 G が全部で m 個の領域から構成されているとする.ここで,G内 のある2つ以上の領域が連結して形成される部分集合を集積領域群の候補と考え、この部分集 合をウィンドウと呼び、Z で表す.このとき、任意のZにおいて、 $i \in Z$ を満たす領域iでの 観測値 y_i は、 $y_i|w_i \sim N(\mu_z, \sigma^2/w_i)$ に従うものと仮定する.また、 $i \in Z^c (= G - Z)$ において は、 $y_i|w_i \sim N(\mu_{z^c}, \sigma^2/w_i)$ に従うとする.なお、これらの観測値は互いに独立であるとする. ここに、 μ_z と μ_{z^c} はそれぞれの領域集合における平均、 σ^2 は領域全体の分散を表す.また、 観測値の地域信頼性や不確実性の度合いを調整するために, 領域 i に関連付けられた重み w_i を 設定する.これにより、個々の観測結果ではなく、各領域内の観測値の全体的な振る舞いを反 映した尺度に基づいて,空間集積性の有無を評価することが可能となる.ここで,wiが大きい ほど,その領域で観測される値の振れ幅が小さくなり、yi は領域 i の真値に近い,信頼性の高 い値と見なされる. wi には、領域 i における観測値の分散の逆数や観測数などを充てることが できる.本稿では、各領域i(i = 1, 2, ..., m)において一組の (u_i, w_i) が利用可能な状況を対象 とする. なお,各領域 i において1つのデータ y_i しか観測されていない場合には,重み $w_i = 1$ を一律に与えるか,あるいは yi の信頼性(重み wi)を外生的に与える方法が考えられる.具体 的には、地域ごとの自然災害リスクやインフラ整備状況などの地理的要因を指標化したり、歴 史的背景といった過去の情報に基づいて重みを設定する方法が挙げられる、さらには、専門家 の意見や評価に基づいてデータの信頼性を外生的に付与するアプローチも考えられる.

さてこのとき, Z に高い観測値からなる集積性(ホットスポットクラスター)が認められるか 否かは,次の仮説検定問題となる.

$$\begin{aligned} H_0 : \mu_{\boldsymbol{z}} &= \mu_{\boldsymbol{z}^c} = \mu_0 & \text{for } {}^{\forall} \boldsymbol{Z} \in \mathcal{Z} \\ H_1 : \mu_{\boldsymbol{z}} &> \mu_{\boldsymbol{z}^c} & \text{for } {}^{\exists} \boldsymbol{Z} \in \mathcal{Z} \end{aligned}$$

ここに, *Z*は *Z*の取り得る全体集合である.

空間スキャン統計量は、 H_0 と H_1 の 2 つのモデルの尤度比で定義される.いま、 H_1 のもと での尤度関数は、(2.1)式で表される.

(2.1)
$$L_{1}(\mathbf{Z}) = \prod_{i \in \mathbf{Z}} \sqrt{\frac{w_{i}}{2\pi\sigma^{2}}} \exp\left(-\frac{w_{i}}{2\sigma^{2}}(y_{i}-\mu_{\mathbf{z}})^{2}\right) \times \prod_{i \notin \mathbf{Z}} \sqrt{\frac{w_{i}}{2\pi\sigma^{2}}} \exp\left(-\frac{w_{i}}{2\sigma^{2}}(y_{i}-\mu_{\mathbf{z}^{c}})^{2}\right) \\ = (\sqrt{\sigma^{2}})^{-m} \exp\left[-\frac{1}{2\sigma^{2}} \left(\sum_{i \in \mathbf{Z}} w_{i}(y_{i}-\mu_{\mathbf{z}})^{2} + \sum_{i \notin \mathbf{Z}} w_{i}(y_{i}-\mu_{\mathbf{z}^{c}})^{2}\right)\right] \prod_{i=1}^{m} \sqrt{\frac{w_{i}}{2\pi}}$$

このとき,重み付き平均 $\mu_{z}, \mu_{z^{c}}$ の最尤推定量は,それぞれ $\hat{\mu}_{z} = \frac{\sum_{i \in \mathbb{Z}} (w_{i}y_{i})}{\sum_{i \in \mathbb{Z}} w_{i}}, \hat{\mu}_{z^{c}} = \frac{\sum_{i \notin \mathbb{Z}} (w_{i}y_{i})}{\sum_{i \notin \mathbb{Z}} w_{i}}$ となる.また,重み付き分散 σ^{2} の最尤推定量を $\hat{\sigma}_{1}^{2}$ とすると, $\hat{\sigma}_{1}^{2} = \frac{\sum_{i \in \mathbb{Z}} w_{i}(y_{i} - \hat{\mu}_{z})^{2} + \sum_{i \notin \mathbb{Z}} w_{i}(y_{i} - \hat{\mu}_{z^{c}})^{2}}{\sum_{i=1}^{m} w_{i}}$ と与えることができる.一方, H_{0} のもとでの尤度関数は, (2.2)式で表される.

(2.2)
$$L_{0} = \prod_{i=1}^{m} \sqrt{\frac{w_{i}}{2\pi\sigma^{2}}} \exp\left(-\frac{w_{i}}{2\sigma^{2}}(y_{i}-\mu_{0})^{2}\right)$$
$$= (\sqrt{\sigma^{2}})^{-m} \exp\left[-\frac{1}{2\sigma^{2}}\sum_{i=1}^{m} w_{i}(y_{i}-\mu_{0})^{2}\right] \prod_{i=1}^{m} \sqrt{\frac{w_{i}}{2\pi}}$$

このとき,重み付き平均 μ_0 の最尤推定量は $\hat{\mu}_0 = \frac{\sum_{i=1}^m (w_i y_i)}{\sum_{i=1}^m w_i}$,重み付き分散 σ^2 の最尤推定量 $\hat{\sigma}_0^2$ は $\hat{\sigma}_0^2 = \frac{\sum_{i=1}^m w_i (y_i - \hat{\mu}_0)^2}{\sum_{i=1}^m w_i}$ でそれぞれ与えられる.

尤度比 $LR = L_1/L_0$ は、最尤推定量を代入することで(2.3)式を導出する.

(2.3)
$$LR(\mathbf{Z}) = \begin{cases} (\hat{\sigma}_1^2 / \hat{\sigma}_0^2)^{-m/2}, & \hat{\mu}_{\mathbf{z}} > \hat{\mu}_{\mathbf{z}^c} \\ 1, & \mathcal{E}の他 \end{cases}$$

通常,計算を簡略化するために,*LR*を対数変換した*LLR* = log*LR*が使用される.なお,*Z*が低い観測値からなる集積性(コールドスポットクラスター)であるか否かについて検討する場合は,*H*₁の不等号の向きを逆にして考えればよく,このときの尤度比は(2.3)式の上段の不等号の向きを逆にしたものが対応する.また,これ以降はホットスポットクラスターを "ホットスポット",コールドスポットクラスターを "コールドスポット" と呼ぶこととする.

なおここに, $\hat{\mu}_{z}, \hat{\mu}_{z^{e}}, \hat{\mu}_{0}$ は, それぞれ $\mu_{z}, \mu_{z^{e}}, \mu_{0}$ の不偏推定量となるが, $\hat{\sigma}_{1}^{2}$ および $\hat{\sigma}_{0}^{2}$ は, それぞれ σ_{1}^{2} と σ_{0}^{2} の不偏推定量とはならない. そのため,実際に *LLR*を算出する際には,不 偏性を考慮した分散を用いる.重み付き分散の不偏推定量にはいくつかの流儀が存在するが, 本稿では SaTScanTM が採用している方法,すなわち重み付き分散の最尤推定量に $\frac{m}{m-1}$ を乗じ たものをそれぞれ $\hat{\sigma}_{1}^{2}, \hat{\sigma}_{0}^{2}$ とする.

つづいて、ウィンドウ Z の全体集合 $Z = \{Z_1, Z_2, ...\}$ の中から最大の対数尤度比をもつ、

(2.4)
$$\hat{\boldsymbol{Z}} = \arg \max_{\boldsymbol{Z} \in \mathcal{Z}} LLR(\boldsymbol{Z})$$

なるウィンドウ \hat{Z} をホットスポットの候補と考える.また, \hat{Z} が統計的に有意なホットス ポットであるかどうかを評価するためには, H_0 の下での $\max_{Z \in \mathbb{Z}} LLR(Z)$ の分布が必要だ が,これを一意に定めることは解析的に困難なため、モンテカルロ法によるシミュレーショ ンを用いて求めた p 値によって有意性を評価するのが通例となっている.具体的には、非パ ラメトリックな検定法である permutation test を利用することで,観測値の分布に影響を受け ることなく有意性を評価する方法が行われる (Kulldorff et al., 2009; Huang et al., 2009).こ れは、領域 i(i = 1, 2, ..., m) における (y_i, w_i) の組み合わせを固定した上で、領域 i の配置を ランダムに並び替えることで生成されるデータセット (y_i^*, w_i^*) に対してウィンドウ Z^* の集 合 Z^* を求め、 $\lambda^* = \max_{Z^* \in \mathbb{Z}^*} LLR(Z^*)$ を算出する.同様の手順を N 回繰り返し、得られた $\Lambda^* = {\lambda_{1,1}^*, \lambda_{2,...,\lambda_N}^*}$ を H_0 の下での最大対数尤度比の分布と見なして、元の観測データセッ ト Gから求めた $LLR(\hat{Z})$ が Λ^* の中でどの程度極端な位置にあるかを確認し、その確率を p 値 とするものである. Λ^* における $LLR(\hat{Z})$ の降順での順位を Rとすると、ホットスポットの候 補 \hat{Z} に対する p 値は、

$$(2.5) p = \frac{n}{1+N}$$

となる.

2.2 領域のスキャン

さてここで、ウィンドウ \hat{Z} をどのように求めるかが問題になる. 一つ一つの Z に対し LLR(Z) を算出しながら \hat{Z} を探すのは非常に手間であるし、そもそも領域数 m が極端に少な い場合を除き、一般的に Z を形成するパターンは膨大な数となるため、これらをすべて調べる ことは不可能である。そこで、効率的に Z のパターンを調べる(これを"スキャンする"とい う)ために、現実的に計算可能な M 個の Z の集合 $Z = \{Z_1, Z_2, ..., Z_M\}, Z_i \subset G$ を構築する ためのスキャン方式がいくつか提案されている。空間スキャン統計量を提唱した Kulldorff 自 身は、円状に Z をスキャンする circular scan 法を提案している。これは、ある発生源を中心に 同心円状に感染範囲が拡大していくような伝染性の疾病などを対象にする場合に高い検出力を 示し,また,アルゴリズムが簡便で計算コストが低いといった長所がある反面,川や道路に沿 うような線状やその他の任意の形状をした空間集積性の検出には適していないことが指摘され ている.

3. 位相的階層構造に基づいたスキャン方式

3.1 エシェロン解析法

栗原(2003)や石岡・栗原(2012)は、エシェロン解析法に基づいてデータの位相的な階層構 造を得て、その構造を利用してスキャンを行うエシェロンスキャン法を提案した。エシェロン 解析法(Myers et al., 1997)は、空間データを構成する各領域の構造を客観的に表現するための 手法である。図1は、3×3の格子状のデータに対するエシェロン解析法の適用例を示してい る。この例では、9つの領域の空間的な位置を表面上の1変量値 xの大小(高低)に基づいて分 類し、エシェロンデンドログラム(これ以降、デンドログラムと呼ぶ)と呼ばれるグラフによっ てデータの位相的な階層構造を視覚化している。デンドログラムから、このデータは5つの階 層に分類されることがわかる。

デンドログラムの任意の階層 kについて, kの上位に他に階層が存在しない場合, $k \in \lceil l - 2$ の階層」とみなす。また、ある階層 lの上位に2つ以上の別の階層が存在する場合、 $l \in \lceil 7 > r$ ウンデーションの階層」とみなす。図1のデンドログラムにおいて、階層 1、階層 2、階層 3 は ピークであり、階層 4 と階層 5 はファウンデーションである。エシェロン解析法の最大の利点 は、空間データの各領域に対し近傍情報が与えられれば、次元を問わずその構造が客観的に階



図 1. 3×3の格子状のデータに対するエシェロン解析法の適用例. (a)3×3の格子データ.
 各領域が1変量値 x を有している.ここに, x(1,1) は座標 (1,1)の値を示している.
 (b)xの値に基づいて,データを同じ位相領域に分類する様子. (c)各位相領域の相対
 関係を, x を縦軸にとって表現したもの.ここに, (1,1) は座標 (1,1)の領域を指す.
 (d)デンドログラムの完成.



図 2.3×3の格子状のデータに対しエシェロンスキャン法を適用する様子.ここに、ウィン ドウ内に許容する最大領域数は4としている.このとき、Zは全部で6つのパターン が存在する.

層化できることである.これは、データの分布や構造を調べるために、ヒストグラム、箱ひげ 図、散布図などを使用するのと同様に、空間データを視覚的に記述するための有用なツールと なる.

3.2 エシェロンスキャン法

エシェロンスキャン法では、デンドログラムのピークとなる階層から順に、その階層に含ま れる領域を上から順に Z に取り込みながら $Z = \{Z_1, Z_2, ...\}$ を構築していく.なおここに、解 析者が予め定めた「ウィンドウ内に許容する最大領域数」に達するか、またはその階層に含まれ るすべての領域をスキャンし終えたとき、スキャンの対象が次の階層に移行する.また、ファ ウンデーションの階層をスキャンする場合には、その上位階層に含まれる領域もすべて含めた 状態で、そのファウンデーションの階層の領域を上から順に Z に含めていく.図2は、図1の デンドログラムから構築される Z の結果を示している.なお、エシェロンスキャン法によって 選ばれる Z は、必ず互いに連結する領域群になる。エシェロン解析法やエシェロンスキャン 法のより詳細なアルゴリズムについては、Kurihara et al. (2020)や栗原・石岡 (2021)を参照さ れたい.

4. 数值実験

4.1 実施手法の概要

本節では、「Normal モデルに基づく空間スキャン統計量」と「エシェロンスキャン法」を組み 合わせた新たな空間集積性検出手法の特性について、数値実験を通じて検証する.具体的に

	従来法	提案法
	Normal モデル	Normal モデル
方式	+	+
	Circular scan 法	エシェロンスキャン法
スキャンの際に	タ箔はの広価桂恕	各領域間の近傍情報,
必要となる情報	召頭域の座棕旧報	各領域の1変量値 <i>x</i>
使用した	$SaTScan^{TM}$ (ver.10.1.3),	D/2ト2独白プロガニノ
ソフトウェア	rsatscan (ver.1.0.7)	Rによる独自ノロクノム

表 1. 数値実験で用いる手法の概要.

は、従来広く用いられている circular scan 法との比較を行い両手法の特徴的な違いを明らかに するとともに、母集団分布が正規分布から逸脱した場合での提案法の頑健性についても確認す る.従来法と提案法の概要を表1に示す.

従来法である circular scan 法は, それぞれ領域 i(i = 1, 2, ..., m) を起点として同心円状にス キャンを行う.言い換えれば, circular scan 法は特定の領域 i を中心として, その周囲の領域 を順に調査していくものと言える.ここで,領域間の距離が重要な役割を果たし, i と距離が 近い領域から順にウィンドウ Z に取り込みながら $Z = \{Z_1, Z_2, ...\}$ を構築していく.領域間 の距離の算出には,各領域の代表点の座標情報を利用する.従来法の適用に際し, Kulldorff et al. (2024)が開発したソフトトウェア SaTScanTM および, SaTScanTM を R 上で実行可能にする rsatscan パッケージ (Kleinman, 2023)を使用する.

提案法では、まずデータに対してデンドログラムを作成するために、各領域間の近傍情報と 各領域の1変量値 x_i (i = 1, 2, ..., m)を準備する必要がある. ホットスポットが比較的大きな 値を持つ x で形成される領域群の場合、それらの各領域はデンドログラムの上位階層に位置す ることになるため、優先的に Z に取り込まれる. このプロセスにより、エシェロンスキャン法 ではある特定の形状に囚われない柔軟な形状をした集積領域群の検出を可能にする. x の定義 において最も単純な方法は、領域 i における観測値 y_i をそのまま用いることが考えられる.

(4.1)
$$x_i = y_i, \ i = 1, 2, \dots, n$$

一方で、エシェロンスキャン法は、近傍情報と1変量値に基づいて領域を階層化していく際 に、高い観測値を持つ領域の近傍にある領域を次々と同じ階層に取り込んでしまい、結果とし て集積性の観点からは不自然な、巨大でいびつなホットスポットが検出される場合がある.こ の問題は、空間スキャン統計量のモデルに依存するものではなく、エシェロンスキャン法自 体に起因するものである.なお、この問題に対処するため、竹村 他 (2021)や Takemura et al. (2022)は AESM (Adjusted echelon scan method)を提案している.その着想に基づいて、本稿 では、解析対象領域全体の平均 $\bar{y} = \frac{1}{m} \sum_{i=1}^{m} y_i$ より低い値を持つ領域 *i* については、 $x_i \in y$ の 最小値に置き換えることを試みる.これにより、 \bar{y} を下回る値を持つような、本来ならばホッ トスポットとして適切でない領域 *i* はデンドログラムの底に移動し、スキャン対象から除外さ れることとなる (図 3).このアプローチは、スキャンされるウィンドウ数の削減にも繋がり、 結果的に計算コストの改善にもなる.

また、重み w_i は各観測値の不確実性を示す指標であり、 w_i が大きいほど、その領域iにおける観測値 y_i の変動が小さくなり、 y_i は真値に近い信頼性の高い値と見なすことができる。この特性を踏まえ、重みの大きい、すなわち信頼性の高い観測値を持つ領域ほどデンドログラムの上位に配置されやすく、逆に重みが小さい観測値を持つ領域ほどデンドログラムの下位に配置されやすくするアプローチを試みる。その一環として、ここでは(4.2)式について検討する。



図 3. (a) 図1の3×3の格子状のデータのデンドログラム. (b) *y*よりも低い値を持つ領域の 値 *x* を *y* の最小値に変換した際に作成されるデンドログラム. (c)変換後のデンドログ ラムに対するスキャン.ここに、ウィンドウ内に許容する最大領域数は4としている.

これにより,例えば複数の高い観測値を持つ領域が存在し,それらの重みが異なる場合に,より信頼性の高い領域がホットスポットして検出されやすくなることが期待される.

(4.2)
$$x_i = \begin{cases} w_i \cdot y_i, & y_i > \bar{y} \\ w_i \cdot \min\{y_1, y_2, \dots, y_m\}, & \mathcal{EO} \end{pmatrix}$$

なお,提案法を適用するためのソフトウェアは存在しないが,デンドログラムの階層構造の 取得には,Rの echelon パッケージ (Ishioka, 2024)が利用できる.

4.2 設計

数値実験の設計においては、先行研究 (Huang et al., 2009)を一部参考にし、以下の方法に 従った.使用するデータは、 10×10 の格子データ(m = 100)とし、そのデータ上に図4のよう な3つの形状からなる真のホットスポット Z^+ を想定する.これら各形状に対して従来法と提 案法を適用し、 Z^+ の形状をどれだけ正確に検出できるかを検証する.

円状を想定した Z^+ として,座標 (3,6)を中心に距離が 2 以内に含まれる領域を設定した. これは従来法でも十分に検出可能と考えられる.また,直線状に伸びた線状の Z^+ と,輪のよ



図 4. 想定した真のホットスポット Z+.

表 2. Z⁺ 内外の各領域の y と w. ケース 1, ケース 2, ケース 3 は従来法と提案法を適用. ケース 4 は(4.2)式に基づく提案法を適用.

		y (乱数	生成)			
	Z ⁺ の形状	$i\in {\bf Z}^+$	$i \notin \mathbf{Z}^+$	重み w	c	η
ケース1	円状	$N(c\sqrt{2},1)$		$w_i = \eta, i \in \mathbf{G}$	1.0	1
			N(0,1)		2.0	1
					3.0	1
ケース2	円状, 線状, 環状	$N(c\sqrt{2},1)$		$w_i = U[1 - \eta, 1 + \eta],$ $i \in \mathbf{G}$	2.0	0.00
			N(0,1)		2.0	0.05
					2.0	0.10
					2.0	0.25
					2.0	0.50
ケース3	円状, 線状, 環状	$N(c\sqrt{2},1)$		$w_i = \eta, i \in \mathbf{Z}^+$ $w_i = 1, i \notin \mathbf{Z}^+$	2.0	1
			N(0,1)		2.0	2
					2.0	5
					2.0	10
					2.0	100
ケース4	円状	$D(c\sqrt{2},1)$	D(0,1)			
		D : 正規分布,ラプラス分布,			2.0	0.00
		ロジスティック分	布,一様分布	$w_i = U[1 - \eta, 1 + \eta],$	2.0	0.05
		$D(1+c\sqrt{2},1)$	D(1,1)	$i\in {f G}$	2.0	0.10
		D: 対数正規分布			2.0	0.25
		$D(1+c\sqrt{2})$	D(1)		2.0	0.50
		D : ポアソ	ン分布			

うな環状の Z^+ を想定した場合についても検証する.これらは従来法ではその形状を正確に捉 えることは困難と予想される.次に、 Z^+ 内外の各領域の $_y$ とwをそれぞれ表 2 のように与 える.

ケース1では、 Z^+ 内外での平均差が、従来法と提案法の違いに及ぼす影響を検証する.ケース2では、ホットスポットの有無にかかわらず、領域全体を対象として、各領域で得られる観測値のばらつき(不確実性)の違いが結果に与える影響を確認する.このとき、重みwは一様乱数によって生成し、具体的には、w = 1.0を中心に幅 0.1 から幅 1.0 までの範囲で生成した場合を検証する.さらに、ケース2では Z^+ の形状を円状、線状、環状に設定した場合の違いについても検証する.

ケース 3 では、ホットスポット内における観測値の不確実性の違いが結果に与える影響を確認する.ここで、 η が大きくなるにつれて、ホットスポット内で観測される yの不確実性が低下することを意味する.なお、2.1節で示した yの従う分布からも明らかなように、 Z^+ の内外

で重みが異なる状況は特殊なケースである点に留意する必要がある.具体例として,都市部に おける感染症の拡大が挙げられる.人口密度の高い都市部では感染症が急速に広がりやすく, クラスター(ホットスポット)が形成されやすい状況が存在する.実際,COVID-19のパンデ ミック時には,都市部での感染拡大が顕著であった.そして都市部には大規模な病院や医療機 関が多く存在し,感染率をより正確に推定するのに必要な情報が集まりやすいと考えられる. また,ケース3においても同様に,**Z**⁺の形状を円状,線状,環状と変化させた場合について 検証する.

ケース4では、正規母集団ではない状況における提案法の頑健性について確認する.ここでは、Huang et al. (2009)の先行研究に倣い、母集団分布として正規分布、ラプラス分布、ロジスティック分布、一様分布、対数正規分布、ポアソン分布を仮定し、これら各々に対して(4.2)式に基づく提案法の検出精度を評価する.このとき、 Z^+ 内外で生成したデータの分布 $D(\mu,\sigma^2)$ は、表2に示す通りである.なお、対数正規分布においては平均を0にすることができないため、 Z^+ 内外において平均を $\mu+1$ に設定している.また、ポアソン分布のパラメータ入は平均と分散の両方を表し、分散が0より大きい値を持つことから、 Z^+ 内では $\lambda=1+c\sqrt{2}$, Z^+ 外では $\lambda=1$ と設定している.さらに、ケース4ではケース2と同様に、領域全体で重みを一様乱数で変動させている.

ここで, circular scan 法で用いる座標情報には領域の中心座標を使用し,エシェロンスキャン法で用いる近傍情報には上下左右の4近傍(ルーク型の隣接関係)を使用する.また,ウィンドウ内に許容する最大領域数は,両手法とも全領域の半数に相当する 50 とした.

4.3 評価指標

各ケースにおいて,従来法と提案法の検出精度は,SensitivityとPPV (positive predictive value)の2つ指標によって評価する.Sensitivity,PPV ともに,空間集積性の検出精度を議論する際に広く用いられている (Huang et al., 2007; Ma et al., 2016; Lee et al., 2021; 竹村 他, 2021).Sensitiviry は, Z^+ 内の領域のうち, \hat{Z} によって捉えることのできた割合のことで,次式で定義される.

(4.3)
$$Sensitivity(\hat{Z}) = \frac{\text{length}(Z^+ \cap Z)}{\text{length}(Z^+)}$$

ここに、length() は領域数である. また、 $0 \leq Sensitivity \leq 1$ であり、1に近いほど「 Z^+ をホットスポットとして検出できた」と判断できる. 一方、*PPV* は、 \hat{Z} 内の領域のうち、 Z^+ を含んでいる割合であり、次式で定義される.

(4.4)
$$PPV(\hat{Z}) = \frac{\text{length}(Z^+ \cap \hat{Z})}{\text{length}(\hat{Z})}$$

Sensitivity と同様に $0 \le PPV \le 1$ となり、1 に近いほど「 Z^+ 以外をホットスポットとして検 出しなかった」と判断できる. さらに、ケース 1、ケース 2、ケース 3 では両手法における対数 尤度比 $LLR(\hat{Z})$ についても確認する. $LLR(\hat{Z})$ が大きいほど、高尤度のホットスポットが検出 できたことになる. なお、今回の数値実験においては、両手法ともに(2.4)式を満たすウィンド ウ \hat{Z} のみを「検出されたホットスポット」と見なし、いわゆる第 2 ホットスポットなどについ ては考慮しない. また、検出された \hat{Z} の形状や $LLR(\hat{Z})$ のみに焦点を当て、有意性について は言及しない.

4.4 結果

表2の各ケースについて,繰り返し数を1000回としてyを乱数によって生成した.図5,図





図 5. 従来法(CS)と提案法(ES1, ES2)に対する Sensitivity, PPV, LLRの結果(ケース 1).



図 6. 従来法(CS)と提案法(ES1, ES2)に対する Sensitivity, PPV, LLRの結果(ケース2).

6,図7,図8に,それぞれ各指標の平均 \pm 3×標準誤差の結果を示す.これらの図において, CSは circular scan法(従来法),ES1とES2はそれぞれ(4.1)式と(4.2)式に基づいたエシェロン スキャン法(提案法)の結果を示している.



図 7. 従来法(CS)と提案法(ES1, ES2)に対する Sensitivity, PPV, LLRの結果(ケース 3).



図 8. 提案法(ES2)に対する Sensitivity, PPV の結果(ケース 4).

ケース 1, ケース 2, ケース 3 の結果から,全体的に提案法は従来法よりも *LLR* の大きい *2*,つまりは高尤度なホットスポットを検出できている.このことは,提案法が形状の制限を 受けず柔軟にウィンドウをスキャンすることに起因していると考えられる.

ケース1の結果に注目すると、従来法は Sensitivity、PPV ともに 0.85 以上を示しており、 このことから真のホットスポット Z^+ を高い精度で捉えることができている. なお、 $c \ge 2$ で は提案法も Sensitivity と PPV の両方で 0.8 以上の高い値を示した. このことから、ホットス ポットが円状である場合、従来法は提案法よりも Z^+ 内外のわずかな平均差をより正確に識別 できることがわかる.

次に、c = 2に固定し、重みを変化させたケース 2 およびケース 3 に焦点を当てる.まず、従 来法については、両ケースともに円状の Z^+ において高い検出精度を示したが、それぞれ線状 の Z^+ では Sensitivity が低下し、環状の Z^+ では PPV が低下する傾向が見られた.一方、提 案法ではすべての形状において高い Sensitivity を維持しており、 Z^+ に含まれる領域をほぼ確 実に捉えていることが確認できる.

また、ケース2において、提案法は全体的に PPV がやや低い値を示した.この原因として、 4.1 節でも述べたように、エシェロンスキャン法の特性上、高い観測値を持つ領域の近傍にあ る領域がスキャン時にウィンドウ内に取り込まれやすいことが挙げられる.しかし、ケース3 のように Z^+ 内の重みが大きくなるにつれて、提案法の PPV が改善する傾向が見られた.な お、重みが極端に大きくなる (ケース3、 $\eta = 100$)場合、従来法の Sensitivity と PPV は提案法 に比べ大きく低下する様子も確認された.さらに、ケース3において Z^+ 内の重みを 100 に設 定した場合、すべての評価指標において提案法が良好な結果を示した.また、提案法における デンドログラムの違い(ES1 と ES2)に着目すると、Sensitivity と LLR にはほとんど差異が見 られなかったが、PPV に関しては ES2 において改善が見られた.

最後に、ケース4の結果について考察する. Sensitivity に関しては、ポアソン分布を除 く5つの連続型分布において、いずれも0.8以上の高い精度が確認された. ポアソン分布の Sensitivity がやや低かった要因としては、他の分布と比較してデータのばらつきが大きく、観 測値が離散値であるために、 Z^+ 内の値の変動が大きくなったことが考えらえる. なお、こ の傾向は Huang et al. (2009)の先行研究における従来法でも確認されている. また、ロジス ティック分布において PPV が特に低かった理由として、ロジスティック分布が正規分布に比 べて裾が広く、 Z^+ 外の領域でも比較的高い値が観測されやすいため、それらの領域がウィン ドウ内に取り込まれたと推察される. 一方、対数正規分布では正規分布よりも PPV が向上し ている. これは分布が右に長い裾を持つため、 Z^+ 外の領域で小さな値が観測されやすくなり、 その結果、 Z^+ 外の領域がウィンドウ内に取り込まれにくくなったためと考えられる.

5. 実データへの応用

前節の数値実験の結果から、データが連続値で与えられる状況において、提案法は、Z⁺内 外の平均に一定の差が存在する場合、その形状にかかわらず、Z⁺を構成する領域を安定的に 取りこぼすことなく検出できることが示唆された.また、母集団が正規分布に従わない場合で あっても、提案法は一定の検出精度を保持できることが確認された.本節では、実データに対 する提案法の適用例を紹介する.

5.1 マース川沿岸の亜鉛濃度のクリギング予測値への適用

5.1.1 データの概要

はじめにクリギングによる推定値に対し提案法を適用する例を紹介する.使用するデータは,





図 9. (a) 元データにおける 155 の観測地点とその亜鉛濃度. (b) 推定する地点を表した格子. (c) 観測地点を Google Earth 上に重ねた図.マース川北部で亜鉛濃度が高くなってい る様子が見て取れる.

R の sp パッケージが提供する meuse である.これは、オランダのマース川沿岸の 155 地点で 観測された複数の重金属濃度の情報を収録しており、今回はこの中から表土の亜鉛濃度(zinc) のデータを扱う.このデータ対して空間補間法の一種であるクリギングを適用し、meuse.grid が提供する 40m 四方の中心座標における 3103 地点の亜鉛濃度を推定する.図9は、155 の観 測地点とその亜鉛濃度、および推定する 3103 地点の格子データを示している.

5.1.2 クリギングによる推定

クリギングは、空間的に相関するデータを補間・推定するための統計的手法である.具体的 には、観測データの空間的な変動パターンをモデル化し、周囲の観測データの加重平均を使



図 10. (a) クリギングによる亜鉛濃度の推定結果の可視化. なお,可視化に際し対数値を元の 単位に戻している. (b) クリギング分散の可視化. 図内の緑点は元データの観測地点 を示す.

表 3. 今回使用したクリギング.

項目	内容
押論バリオグラムモデル	指数型モデル: $\gamma(d_{ij}) = heta_1 + heta_2 \left(1 - \exp\left(-rac{d_{ij}}{ heta_3} ight) ight)$
	<i>d_{ij}</i> は領域 <i>i</i> と領域 <i>j</i> との距離
パラメータの推定法	制限付き最尤推定法
推定されたパラメータ	$\hat{ heta}_1=0, \ \ \hat{ heta}_2=0.622, \ \ \hat{ heta}_3=450.00$
異方性を調整するパラメータ	座標の回転角度:0.469,異方性比:2.605
クリギングのタイプ	通常型クリギング

用して未観測地点での値を推定する.クリギング手法の詳細な説明は間瀬(2010)や瀬谷・堤 (2014)にゆずるとして,ここでは AIC が良好な値を示した推定結果(図10)を採用する.表3 に,今回使用したモデルや手法,推定されたパラメータをまとめている.なお,推定に際し て,亜鉛濃度を対数変換したデータを使用することで,クリギングの精度を向上させている. また,実際のクリギングの適用にはRのgeoRパッケージ(Paulo et al., 2024)を使用した.

このように、クリギングを活用することで、限られた少数のデータから対象領域全体の状況 を推定し、その結果を色分けした図として視覚的に確認することが可能となる.しかしなが ら、その図から受ける印象は、階級区分の設定方法に大きく依存し、例えば濃い色で示された 部分がホットスポットであるかのような誤解を招く恐れがある.今回のように、図 10(a)で示 した 3103 地点の亜鉛濃度の推定値に対し、特に高亜鉛濃度となる領域群の存在の有無を明ら かにすることは、データの未観測領域に跨るホットスポットの挙動に関する知見を得るための 一つのアプローチとして位置付けることができると考える.

また、クリギングによって推定された値の不確実性を評価するために、クリギング分散と呼ばれる指標が用られることがある。クリギング分散は、未知の位置での値の推定誤差の大きさを表すもので、推定値の信頼性を評価する重要な尺度となる。通常、クリギング分散が小さいほど、推定された値の信頼性が高いと判断する。今回求まったクリギング分散を図 10(b)に示す。元々の観測地点が少なかった南西部や、東部から南東部にかけてのエリアでは、クリギン

表 4. 従来法と提案法によって検出されたホットスポットのウィンドウと,そのウィンドウ内 外の重み付き平均,重み付き分散,対数尤度比. p 値については,それぞれスキャン方 式に準じた 3 通りの H₀ の分布から算出している.

						<i>p</i> 值		
スキャン方式	Â	$\hat{\mu}_{\mathbf{z}}$	$\hat{\mu}_{\mathbf{z}^c}$	$\hat{\sigma}_1^2$	$LLR(\hat{\mathbf{Z}})$	$\Lambda^*_{\scriptscriptstyle \mathrm{CS}}$	$\Lambda^*_{\scriptscriptstyle{\mathrm{ES2}}}$	$\Lambda^*_{\scriptscriptstyle{\mathrm{CS}}} \cup \Lambda^*_{\scriptscriptstyle{\mathrm{ES2}}}$
	$\hat{\mathbf{Z}}_1$	1319.67	377.40	69403.06	346.18	0.0001	0.0001	0.00005
従来法	$\hat{\mathbf{Z}}_2$	762.63	376.10	79393.40	137.53	0.0001	0.0041	0.00205
	$\hat{\mathbf{Z}}_3$	1015.17	388.17	81795.20	91.29	0.0125	0.1827	0.09755
	$\hat{\mathbf{Z}}_1$	801.82	285.74	41941.35	1127.61	0.0001	0.0001	0.00005
提案法	$\hat{\mathbf{Z}}_2$	1008.00	387.94	81711.78	92.88	0.0116	0.1669	0.08920
	$\hat{\mathbf{Z}}_3$	152.45	402.02	85327.54	0.000	1.0000	1.0000	1.00000



図 11. 亜鉛濃度の予測値に対するホットスポットの検出結果.

グ分散の値が高くなっていることが見て取れる.

今回,空間スキャン統計量で使用する観測値 y_i (i = 1, 2, ..., 3103) には,図 10(a)の亜鉛濃度の推定値を使用し、重み w_i (i = 1, 2, ..., 3103) には、図 10(b)のクリギング分散の逆数を用いた.このとき、 $\hat{\mu}_0 = 396.18, \hat{\sigma}_0^2 = 86752.57$ となる.

5.1.3 ホットスポットの検出

それぞれ circular scan 法(従来法)とエシェロンスキャン法(提案法)によって検出されたホットスポットを表4と図11に示す.なおここに、ウィンドウ内に許容する最大領域数は、両手法ともに全領域数の約半数に相当する1551領域としている.空間情報については、前節の数値実験のときと同様に、circular scan 法で用いる座標情報には領域の中心座標を使用し、エシェロンスキャン法における各領域の近傍情報は上下左右の4近傍とした.また、デンドログラム作成のためのxには(4.2)式を用いた.さらにここで、(2.4)式において得られた \hat{Z} を第1ホットスポット \hat{Z}_1 とし、 \hat{Z}_1 と領域が重複しない $Z \in Z$ のうちで、 $LLR(\hat{Z}_1)$ に次いで大きい対数 尤度比 $LLR(\hat{Z}_2)$ を有するウィンドウ \hat{Z}_2 を第2ホットスポットと定義した.同様に、 $\hat{Z}_1 \cup \hat{Z}_2$ と領域が重複しない $Z \in Z$ のうちで、 $LLR(\hat{Z}_3)$ を有するウィンドウ \hat{Z}_3 を第3ホットスポットとして定義している.

北西部で検出されたホットスポットに注目すると、従来法ではそれを2つの異なる領域群と して検出したが、提案法では1つの大規模な領域群として検出した.提案法では、マース川沿 岸に沿うように南北に伸びた細長いエリアをホットスポットとして捉えており、マース川北部



図 12.2 つのスキャン方式による H₀の下での最大対数尤度比の分布,ならびに今回検出さ れた各ウィンドウにおける対数尤度比の位置.

沿岸地域と亜鉛濃度との関連性が示唆される結果となった.また,LLRの値から,提案法の 方が高い尤度のホットスポットを検出できている.

5.1.4 有意性の評価

2.2 節で述べた permutation test によって、検出された \hat{Z} の有意性を確認する.ここで、 H_0 の分布として想定する Λ^* は、当然ながらスキャン方式に大きく影響される。通常は、使用す るスキャン方式に基づいて Λ^* を得ることになるが、ここでは circular scan 法とエシェロンス キャン法という異なる 2 つのスキャン方式を扱うことから、(I) circular scan 法に基づく Λ^*_{CS} と、(II) (4.2)式のエシェロンスキャン法に基づく Λ^*_{ES2} をそれぞれ生成することとした。繰り 返し数については、(I)、(II) ともに N = 9999 とした。図 12 は、クリギングによる亜鉛濃度の 予測データから得られた(I) と(II) の分布と、検出された各ウィンドウにおける対数尤度比の位 置を示している。また、(I) と(II) を合わせた $\Lambda^*_{CS} \cup \Lambda^*_{ES2}$ の分布についても確認した(このと き、N = 9999 + 9999 = 19998 となる)。これら各々の分布に対して p 値を算出した結果を表4 に示す。この結果から、北西部で検出されたホットスポットは高度に有意であることがわかる。

5.2 首都圏の市区町村別合計特殊出生率のベイズ推定値への適用

5.2.1 データの概要

つづいて,首都圏¹⁾の市区町村別に集計された合計特殊出生率に対し,提案法を適用する例 を紹介する.合計特殊出生率とは、「15歳から49歳までの女性の年齢別出生率を合計したも の」(厚生労働省,2024)であり、1人の女性が生涯に産む子供の数を示す重要な指標となる.近 年,日本における出生率の低下は大きな社会問題となっており、政府や自治体、民間団体が 様々な対策を講じている.2024年6月に厚生労働省が公表したデータによると、2023年の日 本における合計特殊出生率は1.20であり、これは8年連続で前の年を下回る結果であった.ま た、都道府県別では東京都の合計特殊出生率が最も低く、0.99であった(NHK WEB,2024).

使用するデータとして、e-Stat「人口動態統計特殊報告」における統計表「合計特殊出生率・母の年齢階級別出生率,都道府県・保健所・市区町村別」から、それぞれ平成 20~24年、平成 25~29年、平成 30年~令和4年の3期間のものをダウンロードした(e-Stat, 2024). これらの データは、国勢調査の年を中心とした5年間の日本における日本人データ²⁾を基に、市区町村



図 13. 厚生労働省公表による首都圏市区町村別の合計特殊出生率のベイズ推定値とその標準 偏差を地図上に可視化した結果.(a)平成 20~24年.ここに、神奈川県相模原市の緑 区・中央区・南区については欠損値(白色)となっている.(b)平成 25~29年.(c)平 成 30年~令和4年.(d)3期間での標準偏差.例えば、茨城県鹿嶋市の合計特殊出生 率のベイズ推定値は、平成 20~24年:1.77、平成 25~29年:1.79、平成 30~令和3 年:1.49であり、その標準偏差は 0.169となる.

ごとで集計されたものとなっている.また、人口の少ない地域における出生数の変動を緩和するために、実際のデータは合計特殊出生率をベイズ推定した値が収録されている.図13は、それぞれの期間における合計特殊出生率のベイズ推定値、ならびに3期間でのばらつき(標準偏差)を地域別に可視化した結果である.これらのデータに対し、それぞれ合計特殊出生率の高い地域群(ホットスポット)と、逆に低い地域群(コールドスポット)の検出を試みる.ここで、3期間のベイズ推定値の平均値をyとし、重みwには標準偏差の逆数を用いる.これらの y_i と w_i (i = 1, 2, ..., 373) に基づいて、 $\hat{\mu}_0 = 1.318, \hat{\sigma}_0^2 = 1.973 \times 10^{-2}$ となる.

5.2.2 空間集積性の検出

ホットスポットの検出については、これまでと同じくデンドログラム作成のための 1 変量 値 x に(4.2)式を用い、(2.3)式に基づいた対数尤度比 *LLR* によって統計量を算出する.一方、 コールドスポットの検出に際しては、合計特殊出生率 y_i が低い地域 i ほど、また、その重み w_i が大きいほど、デンドログラムの上位階層に位置させたいと考え、x には次の(5.1)式を用いる こととした.

(5.1)
$$x_{i} = \begin{cases} -(y_{i}/w_{i}), & y_{i} < \bar{y} \\ -(\max\{y_{1}, y_{2}, \dots, y_{m}\}/w_{i}), & \mathcal{FO} \notin \end{cases}$$

また,対数尤度比については(2.3)式の上段の不等号を逆向きにしたものに基づいて算出する. なおここで,デンドログラム作成のための近傍情報には,市区町村の境界線を共有している地 表 5. 空間集積性が認められたウィンドウ **2**に対する,ウィンドウ内外の重み付き平均,重 み付き分散,対数尤度比, p 値.

集積性のタイプ	$\hat{\mu}_{\mathbf{z}}$	$\hat{\mu}_{\mathbf{z}^c}$	$\hat{\sigma}_1^2$	$LLR(\hat{\mathbf{Z}})$	<i>p</i> 値
ホットスポット	1.496	1.297	1.588×10^{-2}	40.478	0.0004
コールドスポット	1.215	1.358	1.563×10^{-2}	43.416	0.0001



図 14. 有意な空間集積性が認められた地域群. それぞれ赤色がホットスポット, 青色がコー ルドスポットを示す.

域を互いに近傍とした.また、スキャンに関して、ウィンドウ内に許容する最大領域数は全地 域の半数に相当する 186 地域とし、モンテカルロ法の繰り返し数は N = 9999 とした.

上記の設定を基にエシェロンスキャン法を適用した結果,それぞれ有意なホットスポットと コールドスポットが一つずつ認められた.それらの詳細を表5に,地図上に可視化した結果を 図14に示す.

今回解析対象とした首都圏の15年にわたる期間中で,安定して高い合計特殊出生率を示し た地域群(ホットスポット)には,主に山梨県内の市町村によって構成される23地域が同定さ れ,山梨県外では唯一,東京都の檜原村だけが含まれた.また,ホットスポットの範囲に着目 すると,檜原村以外では県境で明確に区切られる結果となった.逆に,安定して低い合計特殊 出生率を示した地域群(コールドスポット)は,山梨県を除いた都県にまたがる78市区町で形 成され,特に東京都,埼玉県,千葉県の地域が比較的多く含まれる結果となった.全体で3,000 万人以上の人口を擁し,日本でも有数の人口密集地域として知られる南関東エリア(東京都,埼 玉県,千葉県,神奈川県)のうち,神奈川県でコールドスポットとして認められたのは川崎市 中原区の1地域のみであった.

6. まとめと今後の課題

本稿では、連続値からなる空間データに対し、「Normal モデルの空間スキャン統計量」と「エ シェロンスキャン法」を組み合わせた新たな空間集積性の検出を試みた.また、提案法の有用 性を数値実験で検証した結果、従来法よりも高い尤度を持ち、柔軟な形状からなる集積領域群 を検出できることが示された.特に、集積領域群内の重みが大きい(集積領域群内で観測され るデータが安定していて、ばらつきが少ない)場合は、その傾向が顕著であった.提案法の応 用として、まずクリギングによる予測値に適用し、従来法では捉えられなかったホットスポッ トの傾向を明らかにした.また、合計特殊出生率のベイズ推定値に対する適用では、地図上で 示される地理空間データに対し、ホットスポットとコールドスポットの両方の観点から空間集 積性を検討した.以上の結果から,「離散値をとらない空間データに対し,任意の形状をした 集積領域群を評価できるか」という課題に対応できるツールの一つとして提案法は有益と考え られ,今後はさらなる広範な応用が期待される.なお,今回のようにモデルを用いた推定値に 対して適用する際,第一段階で得られる推定値が後続の集積領域検出に大きく影響する.その ため,第一段階で用いるモデルの選定とその推定精度は,全体の結果を左右する極めて重要な 要素となる.この点を踏まえ,例えば「標高や湿度,風向きといった要因の影響を除去し,純 粋に気温の高い領域群を検出する」といった,観測データの傾向を適切に反映したモデルの構 築が必要である.こうしたモデルを基に提案手法を適用することで,目的とする空間集積性の 評価をより高い精度で行うことが可能になるだろう.

次に今後の課題について述べる.エシェロンスキャン法では,まず始めにデンドログラム作成のために各領域に対し1つの変量値 x を定める必要があるが,今回はこの変量値として(4.2) 式や(5.1)式を用いた.これらの式は,データの平均値と最小値(または最大値)を基に簡便に 算出できるという利点があるが,同時にいくつかの課題も抱えている.特に,平均値と最小値 (または最大値)だけではデータの全体的な分布やばらつきを十分に反映できないため,他の統 計指標や手法を取り入れることが考えられる.さらに,これらの式は重みの影響を強く受ける 可能性があり,その点においても慎重な検討が求められる.デンドログラムを作成する際に, 重みをどのように活用するかについては,適用する状況に応じた最適なアプローチを検討する 必要があるだろう.

最後に、本文中でも触れたようにホットスポット候補となる領域の組み合わせの数は、領域 数に対して指数関数的に増加する.そのため,領域の形や大きさを限定して効率的にスキャン を行う方法が数多く研究されてきた.今回提案した手法においても,真に最大尤度となる領域 群を完全に検出することはできない.この点は,他のスキャン手法を含め,空間スキャン検定 の限界および課題の一つであると言える.この問題に対して, Ishioka et al. (2019)では,二分 決定グラフを用いた手法を導入し、領域の形や大きさを限定せずに、すべてのホットスポット 候補領域の組み合わせを網羅的に算出する試みがなされた。しかし、現行の空間スキャン検定 の枠組みにおいては、尤度最大化を目的とするため、結果として巨大でいびつな形状を持つ ホットスポットが検出されやすく、これが結果の解釈を困難にする要因となる場合がある.こ の問題を回避するため、Tango (2008)のように統計量自体に制約を導入した手法や、Moon et al. (2023)のようにホットスポットの大きさに対する最適化について検討がされている.上述 のように、本提案法にはエシェロンデンドログラムの構築という追加のステップが含まれてお り、このデンドログラムの形状が結果に大きく影響する。この点はデメリットとして捉えられ ることもあるが、デンドログラムの情報を有効活用することで、新たなメリットを生み出せる 可能性もあると考える、今後は、異なるデータセットや条件下での有用性についても検証を行 い、より多角的なアプローチを取り入れながら、さらなる検出精度の向上を図ることが重要と なる.

謝 辞

本研究は JSPS 科研費 JP21K11786, JP24K14856 の助成を受けたものである.

注.

¹⁾ 東京都,神奈川県,千葉県,埼玉県,茨城県,群馬県,栃木県,山梨県の1都7県.ただ し空間的な関係性の観点から東京都島嶼部は除いている.このとき,地域数 *m* は 373 と なる. ²⁾ 厚生労働省によると「日本に住んでいる日本人に係る日本において発生した各事象の 全数」.



- Bhatt, V. and Tiwari, N. (2014). A spatial scan statistic for survival data based on Weibull distributiona, Statistics in Medicine, 33, 1867–1876.
- Costa, M. A., Assunção, R. A. and Kulldorff, M. (2012). Constrained spanning tree algorithms for irregularly-shaped spatial clustering, *Computational Statistics and Data Analysis*, 56, 1771– 1783.
- Duczmal, L. and Assunção, R. A. (2004). A simulated annealing strategy for the detection of arbitrarily shaped spatial clusters, *Computational Statistics and Data Analysis*, 45, 269–286.
- e-Stat (2024). 人口動態統計特殊報告, https://www.e-stat.go.jp/stat-search/files?page=1&toukei=004 50013 (最終アクセス日 2024 年 6 月 14 日).
- Huang, L., Kulldorff, M. and Gregorio, D. (2007). A spatial scan statistic for survival data, *Biometrics*, 63, 109–118.
- Huang, L., Tiwari, R., Zuo, J., Kulldorff, M. and Feuer, E. (2009). Weighted normal spatial scan statistic for heterogeneous population data, *Journal of the American Statistical Association*, 104, 886–898.
- Ishioka, F. (2024). echelon v0.2.0: The Echelon Analysis and the Detection of Spatial Clusters using Echelon Scan Method, https://CRAN.R-project.org/package=echelon (最終アクセス日 2024 年 6 月 14 日).
- 石岡文生, 栗原考次 (2012). Echelon 解析に基づくスキャン法によるホットスポット検出について, 統計 数理, 60(1), 93–108.
- Ishioka, F., Kawahara, J., Mizuta, M., Minato, S. and Kurihara, K. (2019). Evaluation of hotspot cluster detection using spatial scan statistic based on exact counting, *Japanese Journal of Statistics* and Data Science, 2, 241–262.
- 掃部耀平, 竹村祐亮, 石岡文生 (2023). エシェロンスキャン法を用いた重み付き Normal モデルに基づく クラスター検出について, 日本計算機統計学会第 37 回シンポジウム講演論文集, 146–149.
- Kleinman, K. (2023). rsatscan v1.0.7: Tools, Classes, and Methods for Interfacing with SaTScan Stand-Alone Software, https://CRAN.R-project.org/package=rsatscan (最終アクセス日 2024 年 6 月 14 日).
- 厚生労働省 (2024). 人口動態統計特殊報告, https://www.mhlw.go.jp/toukei/list/list58-60.html (最終ア クセス日 2024 年 6 月 14 日).
- Kulldorff, M. (1997). A spatial scan statistic, Communications in Statistics: Theory and Methods, 26, 1481–1496.
- Kulldorff, M. and Nagarwalla, N. (1995). Spatial disease clusters: Detection and inference, Statistics in Medicine, 14, 799–810.
- Kulldorff, M., Huang, L. and Konty, K. (2009). A scan statistic for continuous data based on the normal probability model, *International Journal of Health Geographics*, 8, https://doi.org/ 10.1186/1476-072X-8-58.
- Kulldorff, M., Harvard Medical School and Information Management Services Inc. (2024). SaTScan[™]v10.1.3: Software for the Spatial and Space-Time Scan Statistics, http://www. satscan.org (最終アクセス日 2024 年 6 月 14 日).

栗原考次 (2003). 階層的空間構造を利用したホットスポット検出, 計算機統計学, 15, 171–183. 栗原考次, 石岡文生 (2021). 『エシェロン解析: 階層化して視る時空間データ』, 共立出版, 東京.

- Kurihara, K., Ishioka, F. and Kajinishi, S. (2020). Spatial and temporal clustering based on the echelon scan technique and software analysis, *Japanese Journal of Statistics and Data Science*, 3, 313–332.
- Lee, S., Moon, J. and Jung, I. (2021). Optimizing the maximum reported cluster size in the spatial scan statistic for survival data, *International Journal of Health Geographics*, 20, https://doi.org/ 10.1186/s12942-021-00286-w.
- Ma, Y., Yin, F., Zhang, T., Zhou, X.A. and Li, X. (2016). Selection of the maximum spatial cluster size of the spatial scan statistic by using the maximum clustering set-proportion statistic, *PLoS* ONE, 11(1), https://doi.org/10.1371/journal.pone.0147918.
- 間瀬茂 (2010).『地球統計学とクリギング法:R と geoR によるデータ分析』, オーム社, 東京.
- Moon, J., Kim, M. and Jung, I. (2023). Optimizing the maximum reported cluster size for the multinomial-based spatial scan statistic, *International Journal of Health Geographics*, 22, https://doi.org/10.1186/s12942-023-00353-4.
- Myers, W. L., Patil, G. P. and Joly, K. (1997). Echelon approach to areas of concern in synoptic regional monitoring, *Environmental and Ecological Statistics*, 4, 131–152.
- NHK WEB (2024). 去年の合計特殊出生率 過去最低 厚労省「必要な取り組み加速」, https://www3.nhk. or.jp/news/html/20240606/k10014472291000.html (最終アクセス日 2024 年 6 月 14 日).
- Patil, G. P. and Taillie, C. (2004). Upper level set scan statistic for detecting arbitrarily shaped hotspots, *Environmental and Ecological Statistics*, 11, 183–197.
- Paulo, J. R., Peter, D., Ole, C., Martin, S., Roger, B. and Brian, R. (2024). geoR v1.9-4: Analysis of Geostatistical Data, https://CRAN.R-project.org/package=geoR (最終アクセス日 2024 年 6月14日).
- 瀬谷創, 堤盛人 (2014).『空間統計学:自然科学から人文・社会科学まで』, 朝倉出版, 東京.
- 竹村祐亮, 石岡文生, 栗原考次 (2021). Echelon scan 法による高リスクな空間クラスター検出法の提案, 計算機統計学、**34**, 23–43.
- Takemura, Y., Ishioka, F. and Kurihara, K. (2022). Detection of space-time clusters using a topological hierarchy for geospatial data on COVID-19 in Japan, Japanese Journal of Statistics and Data Science, 5, 279–301.
- Tango, T. (2008). A spatial scan statistic with a restricted likelihood ratio, Japanese Journal of Biometrics, 29, 75–95.
- Tango, T. and Takahashi, K. (2005). A flexibly shaped spatial scan statistic for detecting clusters, International Journal of Health Geographics, 4, https://doi.org/10.1186/1476-072X-4-11.

Spatial Clustering Based on Topological Hierarchical Structures for Continuous Values

Fumio Ishioka

Graduate School of Environmental and Life Science, Okayama University

One of the concerns in spatial data analysis is determining whether the subjects of interest are concentrated in a certain region (i.e., whether spatial clustering exists). In recent years, spatial scan statistics, in which each region of spatial data is scanned based on specific rules and evaluated for spatial clustering by its likelihood, have been widely used in various fields. However, many previous studies have primarily used count data (discrete values) as their objects of interest, and the shapes of the clustering regions have also been limited. In this study, we attempt to answer the question, "Can we evaluate arbitrarily shaped spatial clustering for spatial data that does not take discrete values?" by using echelon analysis. The echelon analysis method classifies spatial data into groups of regions (echelons) composed of the same phase based on the univariate value of each region and the neighborhood information between regions, representing them as a graph with a hierarchical structure. By establishing a method for detecting spatial clustering based on this echelon analysis method, it becomes possible to discuss spatial aggregation for data such as estimates based on a certain prediction model, which are often obtained as continuous values. In this paper, the effectiveness of the proposed method is verified through numerical experiments, demonstrating that the proposed method has a more stable and higher likelihood than the conventional method and can detect spatial clustering with flexible shapes. When the proposed method was applied to the predictions by kriging, it revealed trends of spatial clustering that could not be captured by the conventional method. Furthermore, when applied to Bayesian estimates of the total fertility rate, spatial clusterings were identified in terms of both hot-spot and cold-spot clusters in geospatial data presented on a map.

Key words: Spatial clustering, spatial scan statistic, echelon analysis, echelon scan method.