

Estimadores de Semivariância: Uma Revisão

Semivariance Estimators: A Review

Vinícius Basseto Félix*¹, Oilson Alberto Gonzatto Júnior², Diogo Francisco Rossoni³ e Marcos Jardel Henriques⁴

^{1,2,3,4}Departamento de Estatística, Universidade Estadual de Maringá, Paraná, Brasil

Resumo

A Geoestatística é um ramo da estatística responsável pela incorporação e entendimento das dependências espaciais na modelagem de variáveis georreferenciadas. Na busca pelo melhor modelo ajustado, tem-se o desafio de desenvolver e dominar um ferramental que permita a análise — e quantificação — da variabilidade espacial do fenômeno em estudo, por meio de modelagens específicas. Para isso, é comum fazer uso de medidas de correlação espacial como covariância, correlação e, especialmente, semivariância, uma medida resumo da variabilidade e dependência espacial. É então essencial um bom ajuste do semivariograma, e um estimador adequado para a semivariância é necessário para este ajuste. Uma vez que a maioria dos métodos de estimação em Geoestatística e algoritmos de simulação requerem um modelo teórico ajustado a uma semivariância empírica, objetivou-se expor as construções, deduções e a ideia geral que determina a adequabilidade dos principais estimadores de semivariância a fim de divulgar estimadores ainda pouco utilizados pelos pesquisadores. Portanto, este artigo apresenta uma revisão de oito estimadores de semivariância: o estimador clássico de Matheron, Robusto de Cressie e Hawkins, das Medianas de Cressie, de Pairwise, New-1 (MW1) e New-2 (MW2), das Diferenças de Haslett e o estimador Altamente Robusto de Genton.

Palavras-chave: Geoestatística, Estimador, Semivariância, Robustez, Revisão.

Abstract

Geostatistics is a branch of statistics focusing in understanding and modeling the spatial dependence. There are a lot of tools in Geostatistics to quantify the spatial dependence, such as, covariance, variance and semivariance. The semivariance is a critical tool to this job. This paper presents a review of eight estimators of semivariance: Matheron, Robust of Cressie and Hawkins, the Medians of Cressie, of Pairwise, New-1 (MW1) and New-2 (MW2), Haslett differences, Highly Robust of Genton. Besides, the constructions, deductions and general ideias of these estimators are presented.

Keywords: Geostatistics, Estimator, Semivariance, Robust, Revision.

*Autor para correspondência: felix_prot@hotmail.com

Recebido: 03/03/2016 Revisado: 15/07/2016 Aceito: 03/09/2016

1 Introdução

Etimologicamente, o termo Geoestatística designa o estudo estatístico de fenômenos naturais. De uma maneira mais objetiva, a Geoestatística é um ramo da Estatística Espacial que usa o conceito de funções aleatórias para incorporar a dependência espacial nas análises estatísticas (Rossoni et al., 2014).

Segundo Clark (1979), uma medida tomada em um ponto em \mathbb{R}^n (espaço n -dimensional) guarda relações de dependência com medidas tomadas em pontos adjacentes sugerindo uma estrutura de correlação.

Dessa forma, a Geoestatística pressupõe que uma variável aleatória (v.a), tomada em um ponto qualquer, pode ser expressa pela soma de três componentes (Cressie, 1993):

1. um componente estrutural, correspondente a um valor médio ou a uma tendência;
2. um componente aleatório, espacialmente correlacionado;
3. um ruído aleatório.

Assim, seja $Z(s)$ uma variável aleatória em que s denota uma posição em uma, duas ou mais dimensões ($s \in \mathbb{R}^n$), então

$$Z(s) = \mu(s) + \varepsilon'(s) + \varepsilon''(s), \quad (1)$$

em que

- $\mu(s)$ é uma função determinística que representa o componente estrutural;
- $\varepsilon'(s)$ é um termo estocástico que varia localmente e é espacialmente correlacionado;
- $\varepsilon''(s)$ é um ruído aleatório, não correlacionado.

Com o intuito de compreender e modelar o componente aleatório e espacialmente correlacionado, utiliza-se algumas medidas de correlação espacial, tais quais, covariância, correlação e semivariância.

Dentre estas, a semivariância é uma ferramenta crítica para os estudos de Geoestatística, já que:

1. é uma ferramenta para analisar e quantificar a variabilidade espacial do fenômeno em estudo;
2. a maioria dos métodos de estimação Geoestatística e algoritmos de simulação requerem um modelo teórico ajustado a uma semivariância empírica (Gringarten e Deutsch, 2001).

Vários métodos foram propostos ao longo das décadas para estimar a semivariância. O método mais utilizado, comumente chamado de estimador clássico

de semivariância, foi proposto por Matheron (1962), cuja expressão (Journel e Huijbregts, 1978) é dada por

$$2\hat{\gamma}(h) = \frac{1}{N(h)} \sum_{i=1}^{N(h)} [Z(s_i + h) - Z(s_i)]^2,$$

em que

- $\hat{\gamma}(h)$ é o valor da estimativa da semivariância;
- $Z(s_i)$ é o valor da variável no ponto s_i ;
- $Z(s_i + h)$ é o valor da variável no ponto $s_i + h$;
- $N(h)$ é o número de pares separados pelo vetor h .

Todavia, devido a dificuldade para lidar com fenômenos que apresentavam *outliers*, Cressie e Hawkins (1980) propuseram concomitantemente o estimador robusto e o estimador das medianas para a semivariância. Posteriormente, outros estimadores foram propostos, como o estimador de Pairwise (Srivastava e Parker, 1989; Isaaks e Srivastava, 1989); estimador New-1 e estimador New-2 (Li e Lake, 1994); estimador das diferenças de Haslett (Haslett, 1997); e o estimador altamente robusto de Genton (Genton, 1998).

Em sua dissertação Teixeira (2013) fez uma revisão dos estimadores aqui abordados realizando um estudo de simulação computacional para avaliar suas principais características em uma variedade de cenários. Os resultados de sua pesquisa foram também expostos no trabalho (Teixeira e Scalón, 2013).

A abordagem aqui definida tem sua atenção voltada à origem dos estimadores que, por sua vez, caracteriza a adequação do uso, o que permite ao pesquisador realizar a escolha do estimador mais apropriado ao seu problema.

O objetivo deste artigo é divulgar e ampliar a discussão sobre estimadores robustos, além de tornar mais acessível o entendimento da origem dos estimadores de semivariância e, para isso, tenta apresentar de maneira mais clara as deduções expostas nos trabalhos originais, evidenciando seus principais aspectos e desenvolvendo alguns pontos menos imediatos de suas demonstrações.

2 Estimadores

2.1 Estimador clássico de Matheron

O estimador de semivariância clássico, proposto por Matheron (1962), foi construído a partir das médias dos quadrados dos incrementos do processo, e é determinado com base no método dos momentos, que leva em conta duas definições básicas:

1. Seja U uma variável aleatória qualquer, o r -ésimo momento populacional de U , denotado por μ'_r , é

definido como

$$\mu'_r = \mathbb{E}[U^r],$$

se a esperança existir.

2. Seja (U_1, U_2, \dots, U_n) uma amostra aleatória de uma variável aleatória U qualquer, então o r -ésimo momento amostral centrado em 0, denotado por M'_r , é definido como

$$M'_r = \frac{1}{n} \sum_{i=1}^n U_i^r.$$

A obtenção dos estimadores com base no método dos momentos se dá com determinação de igualdades entre os momentos populacionais e amostrais de uma variável aleatória (Ramachandran e Tsokos, 2009). Para deduzir o estimador clássico para a semivariância considere a amostra aleatória $(U_1, \dots, U_{N(h)})$, tal que a variável aleatória U é da forma $U = [Z(s+h) - Z(s)]^2$.

Note que o primeiro momento populacional de U é dado por

$$\begin{aligned} \mu'_1 &= \mathbb{E}[U] \\ &= \mathbb{E}\{[Z(s+h) - Z(s)]^2\} \\ &= 2\gamma(h), \end{aligned}$$

pois, por definição

$$2\gamma(h) = \text{Var}[Z(s+h) - Z(s)],$$

e, tendo em mente a suposição de estacionariedade da média, tem-se $\mathbb{E}[Z(s+h)] = \mathbb{E}[Z(s)]$, portanto, é possível perceber que

$$\begin{aligned} \text{Var}[Z(s+h) - Z(s)] &= \\ &= \mathbb{E}\{[Z(s+h) - Z(s)]^2\} - \underbrace{\left\{\mathbb{E}[Z(s+h) - Z(s)]\right\}^2}_0 \\ &= \mathbb{E}\{[Z(s+h) - Z(s)]^2\}, \end{aligned}$$

desse modo, tem-se a igualdade

$$2\gamma(h) = \mathbb{E}\{[Z(s+h) - Z(s)]^2\}.$$

Por outro lado, o primeiro momento amostral de U é dado por

$$\begin{aligned} M'_1 &= \frac{1}{N(h)} \sum_{i=1}^{N(h)} U_i \\ &= \frac{1}{N(h)} \sum_{i=1}^{N(h)} [Z(s_i+h) - Z(s_i)]^2. \end{aligned}$$

Agora, igualando as expressões obtidas para o momento populacional μ'_1 e para o momento amostral M'_1 tem-se a igualdade

$$2\hat{\gamma}(h) = \frac{1}{N(h)} \sum_{i=1}^{N(h)} [Z(s_i+h) - Z(s_i)]^2, \quad (2)$$

que estabelece o estimador proposto por Matheron (1962), em que

- $\hat{\gamma}(h)$ é o valor da estimativa da semivariância;
- $Z(s_i)$ é o valor da variável no ponto s_i ;
- $Z(s_i+h)$ é o valor da variável no ponto s_i+h ;
- $N(h)$ é o número de pares separados pelo vetor h .

2.2 Estimador Robusto de Cressie e Hawkins

Uma característica marcante do estimador clássico de Matheron está no fato de que ele é influenciado pela presença de *outliers*, uma vez que suas estimativas são determinadas com base no quadrado da diferença entre duas quantidades, o que pode afetar seriamente seus resultados em situações que apresentam qualquer diferença expressiva.

Nesse contexto, na tentativa de contornar esse problema Cressie e Hawkins (1980) propuseram e avaliaram novos métodos de estimação da semivariância, buscando mais robustez. A ideia inicial de um deles consistia em “remover o quadrado” do estimador proposto por Matheron, a fim de que o mesmo não acentuasse mais a influência negativa dos *outliers*.

Sob a consideração de que $Z(s)$ é um processo gaussiano, tem-se que

$$\frac{Z(s+h) - Z(s)}{\sqrt{2\gamma(h)}} \sim N(0,1),$$

logo,

$$\frac{[Z(s+h) - Z(s)]^2}{2\gamma(h)} \sim \chi_1^2,$$

e conseqüentemente

$$[Z(s+h) - Z(s)]^2 \sim 2\gamma(h) \chi_1^2. \quad (3)$$

A distribuição expressa em (3) é assimétrica e, para obtenção de bons resultados na estimação robusta, buscou-se uma transformação para facilitar o processo, a fim de obter simetria. Assim foram utilizadas as transformações da família $\left\{ [Z(s+h) - Z(s)]^2 \right\}^\lambda$ e após um estudo teórico viu-se que $\lambda = 0,25$ leva a uma distribuição aproximadamente normal, sob a hipótese de que $Z(s)$ seja normalmente distribuído.

A partir desta transformação pode-se considerar a variável aleatória $Y = [Z(s+h) - Z(s)]^{1/2}$. Denotando \bar{Y} como a média aritmética das observações independentes $(Y_1, Y_2, \dots, Y_{N(h)})$, e assumindo a normalidade assintótica das mesmas, sua esperança é expressa por

$$\mathbb{E} \left[\frac{\bar{Y}^4}{2\gamma(h)} \right] = 0,457 + \frac{0,494}{N(h)} + \frac{0,045}{[N(h)]^2}, \quad (4)$$

e pode-se reescrever a equação (4) como

$$\frac{\mathbb{E}[\bar{Y}^4]}{2\gamma(h)} = 0,457 + \frac{0,494}{N(h)} + \frac{0,045}{[N(h)]^2}, \quad (5)$$

expressando \bar{Y}^4 em termos de $Z(s)$ em (5)

$$\frac{\mathbb{E} \left[\left(\frac{1}{N(h)} \sum_{i=1}^{N(h)} \sqrt{\|Z(s_i+h) - Z(s_i)\|} \right)^4 \right]}{2\gamma(h)} = \quad (6)$$

$$= 0,457 + \frac{0,494}{N(h)} + \frac{0,045}{[N(h)]^2},$$

isolando-se $2\gamma(h)$ via equação (6), obtém-se que

$$2\hat{\gamma}(h) = \frac{\left(\frac{1}{N(h)} \sum_{i=1}^{N(h)} \sqrt{\|Z(s_i+h) - Z(s_i)\|} \right)^4}{0,457 + \frac{0,494}{N(h)} + \frac{0,045}{[N(h)]^2}},$$

como o termo $0,045/[N(h)]^2$ é pouco significativo na estimativa, Cressie e Hawkins (1980) sugerem desprezá-lo, e assim um estimador condicionalmente não viesado para $\gamma(h)$ é dado por

$$2\hat{\gamma}(h) = \frac{\left(\frac{1}{N(h)} \sum_{i=1}^{N(h)} \sqrt{\|Z(s_i+h) - Z(s_i)\|} \right)^4}{0,457 + \frac{0,494}{N(h)}},$$

em que

- $\hat{\gamma}(h)$ é o valor da estimativa da semivariância;
- $Z(s_i)$ é o valor da variável no ponto s_i ;
- $Z(s_i+h)$ é o valor da variável no ponto s_i+h ;
- $N(h)$ é o número de pares separados pelo vetor h .

2.3 Estimador das Medianas de Cressie

Cressie e Hawkins (1980) propuseram concomitantemente uma modificação da proposta inicial do estimador robusto, que consiste em substituir a média pela mediana no estimador robusto de Cressie e Hawkins.

A presença de uma observação discrepante pode ainda distorcer a estimativa da média; a mediana, por outro lado, é tolerante a erros grosseiros, visto que, na presença de *outliers* ela se mantém pouco, ou, totalmente inalterada. Portanto, obtém-se o estimador das medianas, análogo ao estimador robusto de Cressie e Hawkins.

Assim, dado um processo gaussiano $Z(s)$, é possível identificar a distribuição qui-quadrado, como visto na expressão (3), logo, considerando Q_θ o quantil empírico de ordem θ e $F_{\chi_1^2}$ o quantil de ordem θ da distribuição χ_1^2 , então

$$2\hat{\gamma}(h) = \frac{\left[Q_\theta \left\{ \sqrt{\|Z(s_i+h) - Z(s_i)\|} : i = 1, \dots, N(h) \right\} \right]^4}{F_{\chi_1^2}}, \quad (7)$$

usando $\theta = 1/2$ na expressão (7) e a ideia da transformação de Y utilizada no estimador robusto de Cressie e Hawkins, tem-se

$$2\hat{\gamma}(h) = \frac{\left[\text{med} \left\{ \sqrt{\|Z(s_i+h) - Z(s_i)\|} : i = 1, \dots, N(h) \right\} \right]^4}{0,457}.$$

em que

- $\hat{\gamma}(h)$ é o valor da estimativa da semivariância;
- $Z(s_i)$ o valor da variável no ponto s_i ;
- $Z(s_i+h)$ o valor da variável no ponto s_i+h ;
- $N(h)$ é o número de pares separados pelo vetor h ;
- $\text{med}\{\cdot\}$ denota a mediana do conjunto $\{\cdot\}$.

2.4 Estimador de Pairwise

O estimador de Pairwise é um estimador de variância relativo. Para reduzir o impacto de alguma eventual discrepância expressiva na diferença $[Z(s+h) - Z(s)]$ do estimador clássico de Matheron (que será acentuada ao ser elevada ao quadrado), sugere-se que o quadrado dessa diferença seja dividido pelo quadrado da média local dos valores.

Como destacado em Teixeira (2013), os autores Srivastava e Parker (1989); Isaaks e Srivastava (1989) afirmam que o estimador de Pairwise contribui para produzir uma melhor visualização da continuidade espacial. Li e

Lake (1994) observam que esse estimador se caracteriza por diminuir o efeito de dados discrepantes e tem como desvantagens a soma dos valores de Z nos pontos s e $s + h$ não poder ser nula e, além disso quando a média se aproxima de zero podem ocorrer semivariogramas com “saltos” na semivariância.

O estimador fica então definido pelo incremento da média dos valores ao quadrado, $[Z(s + h) + Z(s)]/2$, como denominador da diferença dos valores no estimador clássico de Matheron (2), ou seja

$$\hat{\gamma}(h) = \frac{1}{2N(h)} \sum_{i=1}^{N(h)} \left(\frac{Z(s_i + h) - Z(s_i)}{\frac{Z(s_i + h) + Z(s_i)}{2}} \right)^2$$

em que

- $\hat{\gamma}(h)$ é o valor da semivariância;
- $Z(s_i)$ o valor da variável no ponto s_i ;
- $Z(s_i + h)$ o valor da variável no ponto $s_i + h$;
- $N(h)$ é o número de pares separados pelo vetor h .

2.5 Estimador New-1

Os estimadores New-1 e New-2 buscam suprir a imprecisão causada pelo decréscimo do número de pares $N(h)$ observados em decorrência do aumento de h , esta a principal fonte de imprecisão dos estimadores (Yang e Lake, 1988).

Li e Lake (1994) propuseram uma medida integral de semivariância, que é definida por meio do momento de ordem $(d - 1)$ de um semivariograma dentro de uma janela com raio h

$$\gamma(h) = \frac{1}{\int_0^h \zeta^{d-1} d\zeta} \int_0^h \zeta^{d-1} \gamma(\zeta) d\zeta,$$

segundo os autores, esta definição pode ser interpretada como uma média ponderada de $\gamma(\zeta)$ sobre $(0, h)$, com uma função de ponderação dada por

$$\frac{\zeta^{d-1}}{\int_0^h \zeta^{d-1} d\zeta},$$

em que d é a dimensão no espaço euclidiano.

Para deduzir a expressão para o estimador New-1 considere

- $f(x, y)$ uma função de distribuição conjunta das variáveis aleatórias X e Y ;
- $f_x(x), f_y(y)$ as distribuições marginais de X e Y respectivamente;

- $f(x|y)$ e $f(y|x)$ as distribuições condicionais de X dado $Y = y$ e Y dado $X = x$, respectivamente.
- $u(X, Y)$ uma função qualquer das variáveis aleatórias X e Y .

É possível escrever

$$f(y|x) = \frac{f(x, y)}{f_x(x)} \iff f(x, y) = f(y|x)f_x(x), \quad (8)$$

se for desejado determinar a esperança de $u(X, Y)$, tem-se

$$\mathbb{E}[u(X, Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} u(x, y) f(x, y) dx dy, \quad (9)$$

substituindo (8) em (9) é possível perceber que

$$\begin{aligned} \mathbb{E}[u(X, Y)] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} u(x, y) f(y|x) f_x(x) dx dy \\ &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} u(x, y) f(y|x) dy \right] f_x(x) dx \\ &= \int_{-\infty}^{\infty} \mathbb{E}[u(X, Y)|X] f_x(x) dx \\ &= \mathbb{E}\left\{ \mathbb{E}[u(X, Y)|X] \right\}. \end{aligned} \quad (10)$$

Agora, considerando

- $X = Z(s)$;
- $Y = Z(s + \zeta)$;
- $u(X, Y) = [Z(s) - Z(s + \zeta)]^2$.

Tendo em mente que

$$2\gamma(h) = \mathbb{E}\left[(Z(s) - Z(s + \zeta))^2 \right],$$

da nova definição para γ , tem-se

$$\begin{aligned} \gamma(h) &= \frac{\int_0^h \zeta^{d-1} \gamma(\zeta) d\zeta}{\int_0^h \zeta^{d-1} d\zeta} \\ &= \frac{\int_0^h \zeta^{d-1} \mathbb{E}\left[(Z(s) - Z(s + \zeta))^2 \right] d\zeta}{2 \int_0^h \zeta^{d-1} d\zeta} \\ &= \mathbb{E} \left\{ \frac{\int_0^h \zeta^{d-1} \mathbb{E}\left[(Z(s) - Z(s + \zeta))^2 | z(s) \right] d\zeta}{2 \int_0^h \zeta^{d-1} d\zeta} \right\}, \end{aligned}$$

discretizando as integrais, obtém-se

$$\begin{aligned} \gamma(h) &= \\ &= \mathbb{E} \left\{ \frac{\int_0^h \xi^{d-1} \mathbb{E} \left[(Z(s) - Z(s + \xi))^2 | z(s) \right] d\xi}{2 \int_0^h \xi^{d-1} d\xi} \right\} \approx \\ &\approx \mathbb{E} \left\{ \frac{\sum_{q=1}^{\lfloor h/\Delta h \rfloor} q^{d-1} \Delta h^d \mathbb{E} \left[(Z(s) - Z(s + q\Delta h))^2 | z(x) \right]}{2 \sum_{q=1}^{\lfloor h/\Delta h \rfloor} q^{d-1} \Delta h^d} \right\}, \end{aligned}$$

agora, analisando as esperanças separadamente, pelo *Estimador Geral de Máxima Verossimilhança*, para a esperança interna, $\mathbb{E} \left[(Z(s) - Z(s + q\Delta h))^2 | z(s) \right]$, denotada por $\mathbb{E} \{ \circ \}$, tem-se

$$\begin{aligned} \mathbb{E} \left[(Z(s) - Z(s + q\Delta h))^2 | z(s) \right] &\approx \\ &\approx \frac{1}{aq^{d-1}} \sum_{k=1}^{aq^{d-1}} [Z(s_i) - Z(s_i + q\Delta h)_k]^2, \end{aligned}$$

em que

- aq^{d-1} é o número de pontos cuja distância de x_i é menor do que, ou igual a, $q\Delta h$;
- a é um *fator geométrico* que expressa o número de células adicionadas à soma para cada incremento em h ;

Para a esperança externa,

$$\mathbb{E} \left\{ \frac{\sum_{q=1}^{\lfloor h/\Delta h \rfloor} q^{d-1} \mathbb{E} \{ \circ \}}{2 \sum_{q=1}^{\lfloor h/\Delta h \rfloor} q^{d-1}} \right\},$$

denotada por $\mathbb{E} \{ \square \}$, tem-se a seguinte estimativa

$$\mathbb{E} \{ \square \} = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{\sum_{q=1}^{\lfloor h/\Delta h \rfloor} q^{d-1} \mathbb{E} \{ \circ \}}{2 \sum_{q=1}^{\lfloor h/\Delta h \rfloor} q^{d-1}} \right\},$$

organizando as expressões $\mathbb{E} \{ \circ \}$ e $\mathbb{E} \{ \square \}$ obtém-se a

expressão geral para $\hat{\gamma}$, dada por

$$\begin{aligned} \hat{\gamma}_{N_1}(h) &= \\ &= \frac{1}{n} \sum_{i=1}^n \left\{ \frac{\sum_{q=1}^{\lfloor h/\Delta h \rfloor} \frac{q^{d-1}}{a q^{d-1}} \sum_{k=1}^{aq^{d-1}} [Z(s_i) - Z(s_i + q\Delta h)_k]^2}{2 \sum_{q=1}^{\lfloor h/\Delta h \rfloor} q^{d-1}} \right\} = \\ &= \frac{1}{n} \sum_{i=1}^n \left\{ \frac{\sum_{q=1}^{\lfloor h/\Delta h \rfloor} \sum_{k=1}^{aq^{d-1}} [Z(s_i) - Z(s_i + q\Delta h)_k]^2}{2a \sum_{q=1}^{\lfloor h/\Delta h \rfloor} q^{d-1}} \right\}, \end{aligned}$$

ou, em uma notação resumida, o estimador New-1 é definido como

$$\hat{\gamma}_{N_1}(h) = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{1}{2m_i} \sum_{j \in D_{i,h}} [Z(s_i) - Z(s_j)]^2 \right\},$$

em que

- $\hat{\gamma}_{N_1}(h)$ é o valor estimado da semivariância;
- $Z(s_i)$ é o valor da variável no ponto s_i ;
- $Z(s_j)$ é o valor da variável no ponto s_j ;
- n é o número total de observações;
- $j \in D_{i,h} \{j : 0 < |s_i - s_j| \leq h\}$;

Note que o estimador nunca sofre com o decréscimo das observações consideradas quando h aumenta, pois, sua definição leva em conta n e não $N(h)$. Essa propriedade permite que o novo estimador use poucos pontos para obter uma grande precisão.

Para uma estimativa isotrópica, usa-se uma região e não uma única direção, neste caso os dados irregularmente espaçados têm um efeito menor sobre este novo estimador. Além disso, o estimador é mais robusto que outros, uma vez que alguns valores extremos não alteram substancialmente a estimativa.

2.6 Estimador New-2

O estimador New-1 foi um estimador proposto para uma nova definição de semivariância. Já o propósito do estimador New-2 é o desenvolvimento de um estimador para a semivariância teórica, isto é, $\gamma(h) = \mathbb{E} \left\{ [Z(s) - Z(s + h)]^2 \right\}$. A ideia neste caso é determinar uma expressão para $\gamma(h)$, mantendo as características interessantes provenientes da nova definição considerada para o estimador New-1. É possível entender a relação

existente entre $\gamma(h)$ e $\gamma_{N_1}(h)$, uma vez que

$$\begin{aligned} \gamma_{N_1}(h) &= \frac{\int_0^h \xi^{d-1} \gamma(\xi) d\xi}{\int_0^h \xi^{d-1} d\xi} \\ \implies \gamma_{N_1}(h) \int_0^h \xi^{d-1} d\xi &= \int_0^h \xi^{d-1} \gamma(\xi) d\xi \\ \implies \gamma_{N_1}(h) \frac{h^d}{d} &= \int_0^h \xi^{d-1} \gamma(\xi) d\xi \\ \implies \left(\gamma_{N_1}(h) \frac{h^d}{d} \right)' &= \left(\int_0^h \xi^{d-1} \gamma(\xi) d\xi \right)' \\ \implies \gamma'_{N_1}(h) \frac{h^d}{d} + h^{d-1} \gamma_{N_1}(h) &= h^{d-1} \gamma(h) \\ \implies \gamma'_{N_1}(h) \frac{h}{d} + \gamma_{N_1}(h) &= \gamma(h). \end{aligned}$$

Tendo a relação anterior em mente, basta combiná-la com a equação que define o estimador New-1, obtendo, portanto, o estimador New-2 cuja expressão é dada por

$$\hat{\gamma}_{N_2}(h) = \hat{\gamma}_{N_1}(h) + \frac{h}{d} \hat{\gamma}'_{N_1}(h),$$

em que

- $\hat{\gamma}_{N_2}(h)$ é o valor da semivariância estimada pelo New-2;
- $\hat{\gamma}'_{N_1}(h)$ é a derivada de $\hat{\gamma}_{N_1}(h)$ em relação a h ;
- h é o vetor de distâncias;
- d é a dimensão no espaço euclidiano.

Como o estimador New-2 faz seus cálculos com base na derivada de $\gamma_{N_1}(h)$, é sugerido que o método numérico para o cálculo desta derivada seja o *Método da Diferença Centrada*, uma vez que levar em conta toda a informação disponível à volta do ponto central produz estimativas mais precisas para a derivada do que levar em conta a informação disponível em uma única direção arbitrária.

2.7 Estimador das diferenças de Haslett

O estimador das diferenças de *Haslett*, proposto por Haslett (1997), tinha como propósito a utilização em séries temporais com o objetivo principal de reconhecer processos autorregressivos de médias móveis (ARMA).

Seja W_t uma série temporal, o modelo ARMA(p, q) como visto em Cryer e Chan (2008) é definido por parâmetros p e q . O parâmetro p representa o número de

termos autorregressivos em um modelo AR(p), enquanto q se refere ao modelo de médias móveis MA(q).

O modelo AR(p) é definido por

$$W_t = \mu + \alpha_1 W_{t-1} + \dots + \alpha_p W_{t-p} + \varepsilon_t, \quad (11)$$

em que ε_t é um erro aleatório normalmente distribuído, com média 0 e variância σ_ε^2 . Os coeficientes $\alpha_1, \dots, \alpha_p$ são os parâmetros do modelo e μ é uma constante. O modelo AR(p) pode ser representado em termos do operador de retardo B , tal que, $B^j W_t = W_{t-j}$, nesse caso

$$\begin{aligned} W_t &= \mu + \alpha_1 B^1 W_t + \dots + \alpha_p B^p W_t + \varepsilon_t \\ &= \mu + \phi(B) W_t + \varepsilon_t. \end{aligned} \quad (12)$$

O modelo MA(q) é definido por

$$W_t = \mu + \varepsilon_t + \beta_1 \varepsilon_{t-1} + \dots + \beta_q \varepsilon_{t-q}, \quad (13)$$

em que ε_t é erro aleatório normalmente distribuído com média 0 e variância σ_ε^2 . Os coeficientes β_1, \dots, β_q são os parâmetros do modelo e μ é uma constante. Aplicando o operador de retardo em (13), nesse caso para os erros, de modo que $B^j \varepsilon = \varepsilon_{t-j}$, tem-se

$$\begin{aligned} W_t &= \mu + B^0 \varepsilon_t + \beta_1 B^1 \varepsilon_t + \dots + \beta_q B^q \varepsilon_t \\ &= \mu + \theta(B) \varepsilon_t. \end{aligned} \quad (14)$$

Tendo em mente a mesma ideia das equações (11) e (13), pode-se definir o modelo ARMA(p, q), expresso por

$$\begin{aligned} W_t &= \mu + \alpha_1 W_{t-1} + \dots + \alpha_p W_{t-p} + \\ &+ \varepsilon_t + \beta_1 \varepsilon_{t-1} + \dots + \beta_q \varepsilon_{t-q}. \end{aligned}$$

ou ainda, de (12) e (14),

$$\phi(B) W_t = \theta(B) \varepsilon_t + \mu.$$

Para definição do estimador, considerando o processo estocástico W_t , para $t = 1, 2, \dots$ tal que

$$\mathbb{E}[W_{t+h} - W_t] = \delta h,$$

e

$$\text{Var}[W_{t+h} - W_t] = 2\gamma(h),$$

em que $\gamma(h)$ é uma *conditionally negative definite function*, ou seja, uma função condicionalmente negativa definida (Cressie, 1993) e δ é uma constante.

O processo será estacionário para média, no caso em que $\mathbb{E}(W_t) = \mu$ e quando $\delta = 0$, e será estacionário para variância, se $\text{Var}(W_t) = \sigma^2$ for assintoticamente definida, neste caso, desde que

$$\begin{aligned} \text{Var}[W_{t+h} - W_t] &= \text{Var}(W_{t+h}) + \text{Var}(W_t) + \\ &- 2 \text{Cov}(W_{t+h}, W_t), \end{aligned}$$

tem-se que

- $\kappa(h) = \text{Cov}(W_{t+h}, W_t)$;
- $\kappa(0) = \text{Var}(W_{t+h}) = \text{Var}(W_t) = \sigma^2$;
- $\gamma(h) = \kappa(0) - \kappa(h) = \sigma^2 - \kappa(h)$.

Quando o processo é não estacionário para variância, a autocovariância não é definida.

Para um conjunto de dados aparentemente estacionário para média, é usual o uso de estimadores para autocovariância ou semivariograma no lag h , tais como

$$\hat{\mu} = \bar{w},$$

a média amostral;

$$\hat{\kappa}(h) = \frac{1}{n} \sum_{t=1}^{n-h} (w_{t+h} - \bar{w})(w_t - \bar{w}),$$

a autocovariância amostral;

$$\hat{\kappa}_r(h) = \frac{1}{n} \sum_{t=1}^{n-h} (r_{t+h} r_t),$$

a autocovariância dos resíduos (r_t), após aplicação do método dos mínimos quadrados ordinários;

$$\hat{\gamma}(h) = \frac{1}{2(n-h)} \sum_{t=1}^{n-h} (w_{t+h} - w_t)^2,$$

o semivariograma clássico amostral. Assim, considerou-se um estimador alternativo, no caso

$$\tilde{\gamma}(h) = \frac{1}{2(n-h-1)} \sum_{t=1}^{n-h} (d_{ht} - \bar{d}_h)^2, \tag{15}$$

em que $d_{ht} = w_{t+h} - w_t$.

É sabido que a variância da média amostral para dados correlacionados não é σ^2/n e, conseqüentemente $\hat{\kappa}(h)$ é viesado. Mesmo $\tilde{\gamma}(h)$ tendo viés, prefere-se seu uso por distinguir entre não estacionariedade de média e variância na identificação do modelo.

Portanto o estimador das diferenças de Haslett, parte desta ideia e da função de variância, definida por

$$\frac{1}{n-1} \sum_{t=1}^n (w_t - \bar{w})^2,$$

porém para o uso em dados espaciais parte-se do estimador clássico de Matheron (2). Entretanto diferente do estimador clássico de Matheron, não se utiliza a diferença entre as observações $d_h = Z(s+h) - Z(s)$, mas sim a diferença entre d_{hi} , como visto na expressão (15), e a média de todos os pares separados pelo vetor h . Deste modo, define-se o estimador de Haslett pela expressão

$$2\hat{\gamma}(h) = \frac{1}{N(h)-1} \sum_{i=1}^{N(h)} (d_{hi} - \bar{d}_h)^2,$$

em que

- $\hat{\gamma}(h)$ é o valor da estimativa da semivariância;
- $d_{hi} = Z(s_i + h) - Z(s_i)$;
- \bar{d}_h é a média de todos os pares separados pelo vetor h ;
- $N(h)$ é o número de pares separados pelo vetor h .

2.8 Estimador Altamente Robusto de Genton

Segundo Genton (1998), o estimador clássico de Matheron não é robusto contra *outliers* e as modificações propostas por Cressie e Hawkins (1980) não seriam suficientes para uma alta eficiência.

A alta robustez é obtida com base na teoria dos estimadores de escala M , que propõe uma classe de estimadores robustos, obtidos pela minimização da soma de funções, sendo que esta denominação vem do método de máxima verossimilhança (Hampel et al., 1986).

A classe de estimadores de escala M é uma generalização proposta por Huber (1964) para os estimadores de máxima verossimilhança. Esses estimadores são definidos pelo valor $T_n = T_n(X_1, \dots, X_n)$ que estabelece um problema de minimização que consiste em determinar o valor T_n que minimiza a expressão

$$\sum_{i=1}^n \rho(x_i; T_n), \tag{16}$$

ou ainda pela solução da equação implícita

$$\sum_{i=1}^n \psi(x_i; T_n) = 0, \tag{17}$$

em que ρ é uma função real arbitrária definida sobre o cartesiano do espaço paramétrico e o suporte da variável aleatória X , e $\psi(x; \theta) = (\partial/\partial\theta)\rho(x; \theta)$. Note que, no caso em que $\rho(x; \theta) = -\log f(x; \theta)$ tem-se o estimador de máxima verossimilhança comumente utilizado.

Para a dedução desse estimador altamente robusto de Genton considere agora que $N(h)$ é o conjunto de todos os pontos que distam no máximo uma distância h entre si e que N_h denota a cardinalidade de $N(h)$.

(Genton, 1998) supôs um conjunto de observações $V_1(h), \dots, V_{N_h}(h)$ i.i.d., em que $V(h) = Z(s+h) - Z(s)$ é um processo estocástico das diferenças no lag h , que se comporta de acordo com a distribuição de um modelo paramétrico $\{F_\sigma; \sigma > 0\}$, em que $F_\sigma(v) = F(v/\sigma)$, indicando que σ é um parâmetro de escala.

Um estimador de escala M , para σ pode ser expresso como $S_{N_h}(V_1(h), \dots, V_{N_h}(h))$ de σ e é definido implicitamente, como visto em (17), pela equação

$$\sum_{i=1}^{N_h} \psi(V_i(h)/S_{N_h}) = 0,$$

e corresponde assintoticamente a estatística funcional S definida por

$$\int \psi(v/S(F)) dF(v) = 0,$$

em que ψ é real, simétrico e uma função suficientemente regular (Hampel et al., 1986).

A função de influência do estimador M de escala S numa distribuição F é expressa por (Hampel et al., 1986)

$$IF(v, S, F) = \frac{\psi(v/S(F))S^2(F)}{\int \psi'(v/S(F)) dF(v)}.$$

A importância desta função consiste em sua interpretação, que retrata o efeito do estimador de contaminação infinitesimal num ponto v , assim algo relevante é denotado no valor de sensibilidade do erro bruto da função, definido por

$$\gamma^* = \sup_v \{IF(v, S, F)\}.$$

Esta medida permite a mensuração da pior influência que uma pequena quantidade contaminada pode ter sobre o estimador.

Outra medida relevante é o ponto de ruptura do estimador de escala, este indica quantos pontos precisam ser trocados para o estimador tender para infinito, ou para zero, para estimadores M como mostrado por Huber (1996), este é dado por

$$\varepsilon^* = \min \left\{ \frac{-\psi(0)}{\chi(+\infty) - \chi(0)}, \frac{\psi(+\infty)}{\psi(+\infty) - \psi(0)} \right\} \leq \frac{1}{2}.$$

Ao escolher-se

$$\psi(v) = \|v\|^q - \int \|v\|^q dF(v), \quad q > 0,$$

tem-se os estimadores de escala L^q (Genton e Rousseeuw, 1995), que são nunca delimitados, consistentes e normalmente distribuídos assintoticamente, para sua função de influência associada para qualquer valor $q > 0$, ou seja $\gamma^* = \infty$, ou ainda $\varepsilon^* = 0\%$.

Percebe-se que os estimadores clássico de Matheron e de Cressie e Hawkins correspondem, respectivamente, a L^2 e $L^{1/2}$, assim ambos não se enquadram no conceito de robustez, pelos critérios de função de influência e ponto de ruptura.

Assim, buscou-se um estimador que atende aos critérios anteriores. Para a estimação robusta, usualmente utilizava-se o *median absolute deviation* (MAD), dado por

$$MAD_n = b \operatorname{med}_i \{ \|s_i - \operatorname{med}_j s_j\| \},$$

em que b é o fator de consistência, outro estimador também usado para tal utilidade foi o S_n definido pela expressão

$$S_n = b \operatorname{med}_i \{ \operatorname{med}_j \|s_i - s_j\| \}.$$

Entretanto ambos tinham a desvantagem de suas funções terem descontinuidades, diferente do estimador de Hodges e Lehmann (1963), dado por

$$\operatorname{med} \left\{ \frac{s_i + s_j}{2}; i < j \right\},$$

que pode ser visto como versão suave da mediana.

Além deste, Shamos (1977) e Bickel e Lehmann (2012) mencionaram um análogo

$$\operatorname{med} \{ \|s_i - s_j\|; i < j \}, \tag{18}$$

entretanto busca-se um estimador com ponto de ruptura de 50%, isto foi obtido facilmente ao mudar a estatística de ordem do estimador da expressão (18)

$$Q_{N(h)} = b \left\{ \|s_i(h) - s_j(h)\| : i < j \right\}_{(k)},$$

em que

$$k = \left(\frac{\lfloor N_h/2 \rfloor + 1}{2} \right),$$

em que $\lfloor N_h/2 \rfloor$ denota a parte inteira de $N_h/2$.

Este estimador compartilha a propriedade de S_n , isto é, o bom funcionamento para distribuições assimétricas. Além disso a função é suave e sua eficiência para distribuições gaussianas é muito alta, em torno de 82% (Rousseeuw e Croux, 1993).

Assim em Genton (1998), definiu-se o fator de consistência para 2,2191, portanto

$$Q_{N(h)} = 2,2191 \left\{ \|V_i(h) - V_j(h)\| : i < j \right\}_{(k)},$$

em que

- $V_i(h) = Z(s_i + h) - Z(s_i)$;
- $V_j(h) = Z(s_j + h) - Z(s_j)$;
- 2,2191 é o fator de consistência utilizado para uma distribuição gaussiana.

Logo, usa-se as diferenças absolutas entre as diferenças, para o k -ésimo quantil, este valor então é multiplicado pelo fator de consistência (Genton, 1998).

Assim, tem-se que o estimador altamente robusto de Genton é definido como

$$2\hat{\gamma}(h) = (Q_{N(h)})^2.$$

em que

- $\hat{\gamma}(h)$ é o valor da estimativa da semivariância;
- $Q_{N(h)}$ é estimador com ponto de ruptura de 50%.

3 Códigos disponíveis atualmente

Teixeira e Scalon (2013) implementaram, na linguagem R e sob a consideração de um *grid* regular, todos os estimadores aqui descritos. Os códigos podem ser consultados em seu trabalho (Teixeira, 2013), entretanto, a maioria dos estimadores não está implementada e disponibilizada nos *softwares* de análise geoestatística usuais, de modo que a utilização dos mesmos está condicionada ao interesse do pesquisador em implementá-los manualmente.

Na Tabela 1 estão listados alguns *softwares* estatísticos que atualmente permitem o uso de alguns dos estimadores abordados neste estudo.

Tabela 1: Implementações disponíveis dos estimadores abordados neste estudo.

Software (Estimador)	Referência
ArcGIS (2.1)	Johnston et al. (2001)
GS+ (2.1)	Robertson (2008)
MATLAB (2.1; 2.2)	Trauth et al. (2007)
R [fractaldim] (2.8)	Sevcikova et al. (2014)
R [geoR] (2.1; 2.2)	Ribeiro Jr e Diggle (2016)
R [georob] (2.1; 2.2; 2.8)	Papritz (2016)
R [gstat] (2.1; 2.2)	Pebesma et al. (2004)
SAS [VARIOGRAM] (2.1; 2.2)	SAS Institute Inc. (2015)
Surfer (2.1)	Barnes (2003)
Variowin (2.1)	Pannatier (1996)
VESPER (2.1)	Whelan et al. (2002)

Note que todos os *softwares* expostos na Tabela 1 disponibilizam o uso do estimador clássico de Matheron, entretanto somente os *softwares* MATLAB, SAS e R possuem o estimador de Cressie e Hawkins, e somente o R, nos pacotes *fractaldim* e *georob*, apresenta o estimador altamente robusto de Genton.

4 Considerações Finais

Diversos estimadores de semivariância foram propostos com o intuito de contornar uma característica negativa do estimador clássico de Matheron. Cada um deles se embasa em ideias distintas para sua construção, mas todos eles compartilham o objetivo de alcançar um alto grau de robustez.

O uso de um estimador de semivariância diferente do estimador clássico de Matheron e, talvez o de Cressie e Hawkins, é incomum, mesmo em situações em que a aplicação dos demais seja a mais adequada. Essa ocorrência pode ter diversas causas, como o desconhecimento da teoria, ou mesmo a falta de acesso a suas implementações nos *softwares* estatísticos.

Nesse sentido, além de conseguir divulgar alguns estimadores pouco abordados na literatura, expondo suas construções, deduções e, quando possível, o acesso às implementações disponíveis, este estudo encoraja suas utilizações à uma gama de profissionais que precisam fazer uso dessas ferramentas em seu dia-a-dia.

Referências

- Barnes, R. (2003). *Variogram tutorial*. Golden Software, Inc., URL <http://www.goldensoftware.com/variogramTutorial.pdf>.
- Bickel, P. J., Lehmann, E. L. (2012). Descriptive statistics for nonparametric models I: Introduction. Em: Rojo, J. (ed) *Selected Works of E. L. Lehmann*, Springer, pp. 465–471.
- Clark, I. (1979). *Practical geostatistics*. Applied Science, London.
- Cressie, N., Hawkins, D. M. (1980). Robust estimation of the variogram: I. *Mathematical Geology*, 12(2), 115–125.
- Cressie, N. A. C. (1993). *Statistics for Spatial Data*. Wiley Series in Probability and Statistics, New York, USA.
- Cryer, J., Chan, K. (2008). *Time Series Analysis: With Applications in R*. Springer Texts in Statistics, Springer.
- Genton, M. G. (1998). Highly robust variogram estimation. *Mathematical Geology*, 30(2), 213–221.
- Genton, M. G., Rousseeuw, P. J. (1995). The change-of-variance function of M-estimators of scale under general contamination. *Journal of computational and applied mathematics*, 64(1), 69–80.
- Gringarten, E., Deutsch, C. V. (2001). Teacher's aide variogram interpretation and modeling. *Mathematical Geology*, 33(4), 507–534.
- Hampel, F. R., Ronchetti, E. M., Rousseeuw, P. J., Stahel, W. A. (1986). *Robust statistics: The approach based on influence functions*. John Wiley & Sons, Inc, New York.
- Haslett, J. (1997). On the sample variogram and the sample autocovariance for non-stationary time series. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 46(4), 475–484.
- Hodges, J. L., Lehmann, E. L. (1963). Estimates of location based on rank tests. *The Annals of Mathematical Statistics*, 34(2), 598–611.
- Huber, P. J. (1964). Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, 35(1), 73–101.

- Huber, P. J. (1996). *Robust Statistical Procedures*. SIAM: Society for Industrial and Applied Mathematics, Philadelphia, PA, USA.
- Isaaks, E. H., Srivastava, R. M. (1989). *An Introduction to Applied Geostatistics*, vol 21. Oxford University Press, Inc., New York.
- Johnston, K., Ver Hoef, J. M., Krivoruchko, K., Lucas, N. (2001). *Using ArcGIS geostatistical analyst*. ESRI, New York.
- Journal, A. G., Huijbregts, C. J. (1978). *Mining geostatistics*. Academic press.
- Li, D., Lake, L. W. (1994). A moving window semivariogram estimator. *Water Resources Research*, 30(5), 1479–1489.
- Matheron, G. (1962). *Traité de géostatistique appliquée*, vol 14. Editions Technip, Paris.
- Pannatier, Y. (1996). *VARIOWIN: Software for spatial data analysis in 2D*. Springer, New York.
- Papritz, A. (2016). *georob: Robust Geostatistical Analysis of Spatial Data*. URL <https://CRAN.R-project.org/package=georob>, R package version 0.2-3.
- Pebesma, E. J., Graeler, B., Pebesma, M. E. (2004). Multi-variable geostatistics in S: the gstat package. *Computers & Geosciences*, 30, 683–691.
- Ramachandran, K. M., Tsokos, C. P. (2009). *Mathematical Statistics with Applications*. Academic Press: Elsevier, San Diego, California, USA.
- Ribeiro Jr, P. J., Diggle, P. J. (2016). *geoR: Analysis of Geostatistical Data*. URL <https://CRAN.R-project.org/package=geoR>, r package version 1.7-5.2.
- Robertson, G. P. (2008). *GS+: Geostatistics for the environmental sciences*. URL <https://www.gammasdesign.com/>, gamma Design Software, Plainwell.
- Rossoni, D. F., de Lima, R. R., de Oliveira, M. S. (2014). Proposta e validação de testes bootstrap para detecção de anisotropia em fenômenos espaciais contínuos. *Revista da Estatística da Universidade Federal de Ouro Preto*, 3(2), 210–227.
- Rousseeuw, P. J., Croux, C. (1993). Alternatives to the median absolute deviation. *Journal of the American Statistical Association*, 88(424), 1273–1283.
- SAS Institute Inc. (2015). The VARIOGRAM procedure. Em: *SAS/STAT(R) 14.1 User Guide*, SAS Institute Inc., Cary, NC, USA, pp. 9846–9969.
- Sevcikova, H., Percival, D., Gneiting, T. (2014). *fractaldim: Estimation of fractal dimensions*. URL <https://CRAN.R-project.org/package=fractaldim>, R package version 0.8-4.
- Shamos, M. I. (1977). Geometry and statistics: Problems at the interface. Em: *Algorithms and Complexity: New Directions and Recent Results*, Citeseer, Academic Press, Inc., New York, pp. 255–280.
- Srivastava, R. M., Parker, H. M. (1989). Robust measures of spatial continuity. Em: *Geoestatics, on the 3rd. Geostatistical Congress*, Holland, Armstrong, pp. 295–308.
- Teixeira, M. B. R. (2013). Comparação entre estimadores de semivariância. Dissertação (mestrado), UFLA, Lavras, MG.
- Teixeira, M. B. R., Scaloni, J. D. (2013). Comparação entre estimadores de semivariância. *Revista Brasileira de Biometria*, 31(2), 248–269.
- Trauth, M. H., Gebbers, R., Marwan, N., Sillmann, E. (2007). *MATLAB Recipes for Earth Sciences*. Springer, New York.
- Whelan, B. M., McBratney, A. B., Minasny, B. (2002). Vesper 1.5 — spatial prediction software for precision agriculture. Em: Robert, P. C., Rust, R. H., Larson, W. E. (Eds) *Precision Agriculture, Proceedings of the 6th International Conference on Precision Agriculture*, ASA/CSSA/SSSA, Madison, WI, USA.
- Yang, A. P., Lake, L. W. (1988). The accuracy of autocorrelation estimates. *In Situ*, 12(4), 227–274.