

Universidade Estadual de Maringá

Programa de Iniciação Científica – PIC

Departamento de Estatística

Orientador: Prof. Dr. Diogo Francisco Rossoni

Acadêmico: André Felipe Berdusco Menezes

Preditores Geoestatísticos: revisão e aplicação

Maringá, 31 de Julho de 2017

Universidade Estadual de Maringá

Programa de Iniciação Científica – PIC

Departamento de Estatística

Orientador: Prof. Dr. Diogo Francisco Rossoni

Acadêmico: André Felipe Berdusco Menezes

Preditores Geoestatísticos: revisão e aplicação

Relatório contendo os resultados finais do projeto de iniciação científica vinculado ao Programa PIC-UEM.

Maringá, 31 de Julho de 2017

Resumo

Realizar predições de valores desconhecidos de fenômenos compõe uma das principais atividades científicas de muitos pesquisadores. Tendo em vista que o fenômeno em estudo é espacialmente correlacionado, a metodologia geoestatística proporciona distintos preditores. Denominados krigagem, esses preditores usam a vizinhança amostrada e incorporam a estrutura de dependência espacial, além de possuírem as características de não tendenciosidade e variância mínima. Podemos definir o preditor de krigagem como uma média ponderada, sendo os pesos atribuídos conforme o modelo de dependência espacial e o tipo de krigagem empregada. Dessa forma, objetiva-se neste estudo compreender e descrever de modo detalhado os seguintes preditores geoestatísticos: krigagem da média, krigagem simples, krigagem ordinária, krigagem universal e krigagem indicatriz.

Palavras-chave: Geoestatística, Krigagem, Revisão.

Conteúdo

1	Introdução	2
2	Objetivos	3
2.1	Objetivo Geral	3
2.2	Objetivos Específicos	3
3	Desenvolvimento	4
3.1	Contextualização	4
3.2	Materiais e Métodos	5
3.2.1	Conjunto de dados	5
3.2.2	Krigagem da Média	7
3.2.3	Krigagem Simples	9
3.2.4	Krigagem Ordinária	11
3.2.5	Krigagem Indicatriz	13
3.2.6	Krigagem Universal	14
3.3	Resultados e Discussões	15
3.3.1	Krigagem da Média no R	16
3.3.2	Krigagem Simples no R	17
3.3.3	Krigagem Ordinária no R	17
3.3.4	Krigagem Indicatriz no R	18
3.3.5	Krigagem Universal no R	20
4	Conclusão	23

1 Introdução

Neste projeto terá enfoque o estudo da metodologia geoestatística, em específico os tipos de preditores lineares, uma vez que a predição está presente em diversas áreas do conhecimento. Entretanto em muitos estudos, os pesquisadores não fazem uso da informação espacial dos dados, sendo ela necessária para uma melhor predição.

Matheron (1963), fundamentou os conceitos teóricos da geoestatística em processos estocástico, isto é, nas análises considera-se que o fenômeno em estudo é oriundo de um processo estocástico. A compreensão da variabilidade espacial do processo estocástico ocorre por meio de uma modelagem probabilística. Esta modelagem é realizada utilizando-se de medidas de associação próprias da geoestatística, tais como semivariância e covariância.

Todavia, na maioria das vezes, o interesse não se restringe em obter apenas o modelo de dependência espacial, visando também prever/interpolar valores em pontos não amostrados. Existem diversos tipos de interpoladores na literatura, como por exemplo: poligonal, inverso do quadrado das distâncias, triangulação e médias locais. Entretanto estes não consideram a variabilidade espacial do fenômeno. Em contrapartida os preditores geoestatísticos, conhecidos como krigagem e definidos inicialmente por Matheron (1963), fazem uso do modelo de dependência espacial para prever, seja o valor esperado ou a função de distribuição acumulada de variáveis aleatórias do processo estocástico.

De modo geral, podemos afirmar que o preditor de krigagem é uma média ponderada, sendo os pesos atribuídos conforme o modelo ajustado no semivariograma e o tipo de krigagem empregada. Especificamente, neste projeto serão estudados de forma minuciosa os seguintes tipos de krigagem: krigagem da média, krigagem simples, krigagem ordinária, krigagem indicatriz e krigagem universal.

2 Objetivos

2.1 Objetivo Geral

Realizar uma revisão detalhada dos principais tipos de Krigagem, uma vez que este método é considerado o melhor preditor linear não viciado e de variância mínima quando existem evidências de dependência espacial.

2.2 Objetivos Específicos

- Fazer uma revisão sobre krigagem da média, krigagem simples, krigagem ordinária, krigagem indicatriz e krigagem universal.
- Verificar no ambiente R quais tipos de krigagem já estão devidamente implementadas em pacotes distintos;
- Propor algoritmos em linguagem R para as possíveis krigagens que ainda não possuem implementação computacional;

3 Desenvolvimento

3.1 Contextualização

Este projeto é uma continuidade do projeto de iniciação científica 3481/2015, o qual teve por objetivo, propor um método de sub-amostras para o cálculo da semivariância, a qual é fundamental para o ajuste e compreensão de um modelo espacial. Dessa forma, alguns conceitos discutidos de forma detalhada no projeto anterior serão reproduzidos de forma sucinta nesta seção.

Conforme define Matheron (1963), em geoestatística os dados são provenientes de um processo estocástico, isto é, para cada ponto x_i amostrado tem-se uma variável aleatória Z distinta. Formalmente temos:

$$\{Z(x_i) : x \in D \subset \mathbb{R}^p\} \quad (1)$$

Sendo:

- Z a variável aleatória que varia continuamente em D ;
- x a posição da variável, considerada fixa;
- D a região em estudo;
- \mathbb{R}^p o espaço p -dimensional ($p = 1, 2, 3$ ou 4)

Devido a limitações de ordem prática, na maioria das análises geoestatística existe somente uma realização do processo estocástico, o que torna impossível realizar inferência sobre este processo. Dessa forma é necessário que algum tipo de estacionaridade, adequada com o problema em estudo, seja assumida de maneira a possibilitar estimação de ao menos os dois primeiros momentos do processo estocástico, isto é média, covariância e/ou semivariância.

Para compreender a variabilidade do processo estocástico, presume-se que o valor de um ponto no espaço está relacionado, de alguma forma, com valores de pontos situados a certa distância, sendo provável supor que a influência é tanto maior quanto menor for à distância entre eles.

Várias medidas prestam para descrever a variabilidade, tais como covariância, correlograma e semivariância. Em geoestatística, usualmente utiliza-se a semivariância, uma medida de dissimilaridade, isto é seu valor cresce conforme a associação entre as variáveis. O gráfico da semivariância em função da distância é denominado semivariograma.

A modelagem probabilística presente na geoestatística consiste no ajuste de algum modelo teórico ao semivariograma empírico (CRESSIE, 1985). O modelo que melhor se ajusta aos pontos representa a magnitude, o alcance e a intensidade da variabilidade espacial do processo estocástico.

Embora a modelagem seja fundamental para descrição e entendimento do fenômeno, na atividade científica o interesse não se limita em obter apenas o modelo de dependência espacial, desejando-se também prever valores em pontos não amostrados. O pesquisador pode estar interessado em um ou mais pontos específicos ou ainda obter uma malha de pontos interpolados.

A técnica adotada na geoestatística para predição é intitulada krigagem. Ela faz uso da vizinhança amostral implementando o modelo de dependência espacial do processo, reduzindo assim a incerteza nas predições. Associada a krigagem esta

o acrônimo **BLUP** (*Best Linear Unbiased Predictor*). Conforme esclarece Santos (2010), o termo “Linear” indica que as predições são combinações lineares dos dados observados; “Unbiased” significa que o valor esperado dos resíduos é zero; “Best” porque visa minimizar a variância dos resíduos e, “Predictor” por estimar variáveis aleatórias.

Nas seções posteriores abordaremos a krigagem da média, krigagem simples, krigagem ordinária, krigagem indicatriz e krigagem universal.

3.2 Materiais e Métodos

3.2.1 Conjunto de dados

Investigaremos a performance dos diferentes métodos de krigagem analisando o banco de dados **meuse**, junto ao software R. O conjunto de dados **meuse** esta disponível na biblioteca **gstat** (PEBESMA, 2007), ele contém medições de concentrações em mg/kg de diferentes elementos químicos extraídos 164 localizações do rio Meuse na Holanda. Em específico iremos analisar o comportamento do chumbo (Pb).

Inicialmente foi realizada uma análise descritiva da variável chumbo, apresentamos na Figura 1 a distribuição espacial dos dados coletados. É possível verificar alto índice de chumbo, isto é, acima de 300 mg/kg no litoral, enquanto que índices menores são encontrados no centro. Em relação a distribuição da variável verifica-se uma assimetria a direita com cauda pesada, indicando fuga da normalidade.

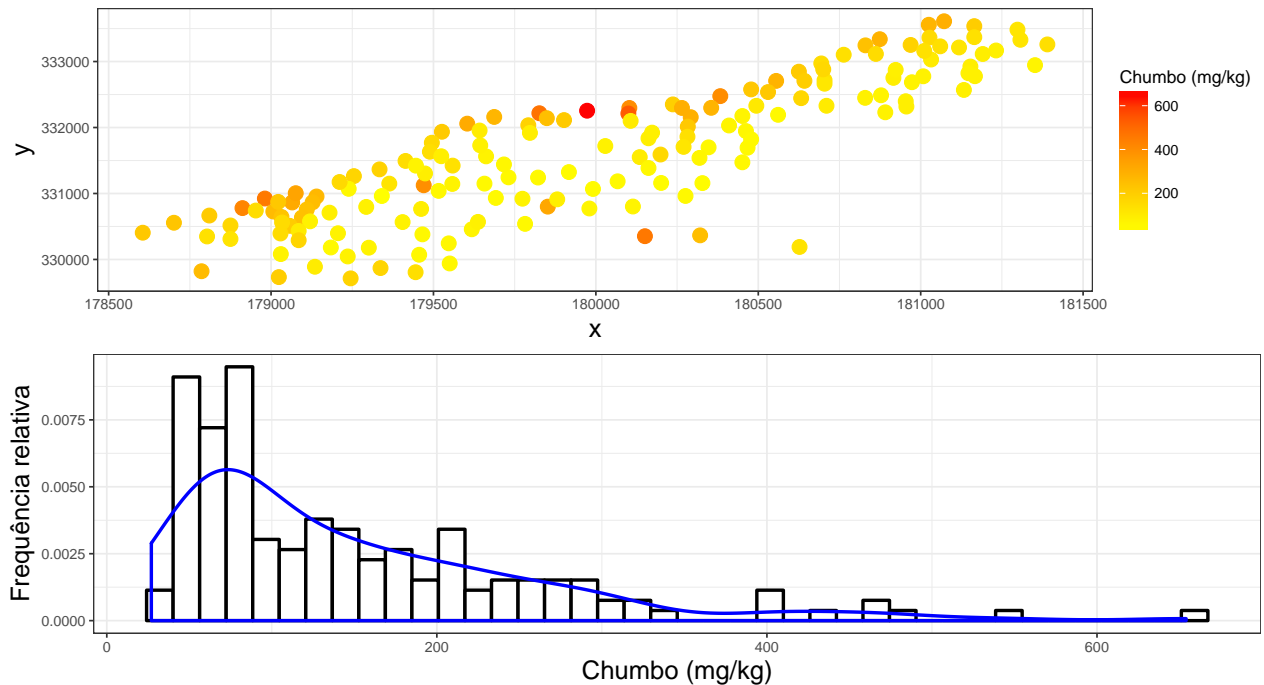


Figura 1: Comportamento do chumbo (mg/kg)

Na análise geoestatística tivemos o interesse de identificar e quantificar o modelo de dependência espacial do chumbo, uma vez que somente a partir dele é

possível realizar predições de pontos não observados. Devemos então, encontrar o semivariograma empírico e ajustá-lo a um modelo teórico.

Nesta etapa as análises foram conduzidas com auxílio do pacote `gstat` (PEBESMA, 2007) em específico das funções `variogram` e `fit.variogram`. Os modelos teóricos considerados foram exponencial, esférico e gaussiano. A Figura 2 apresenta o semivariograma empírico e os modelos ajustados.

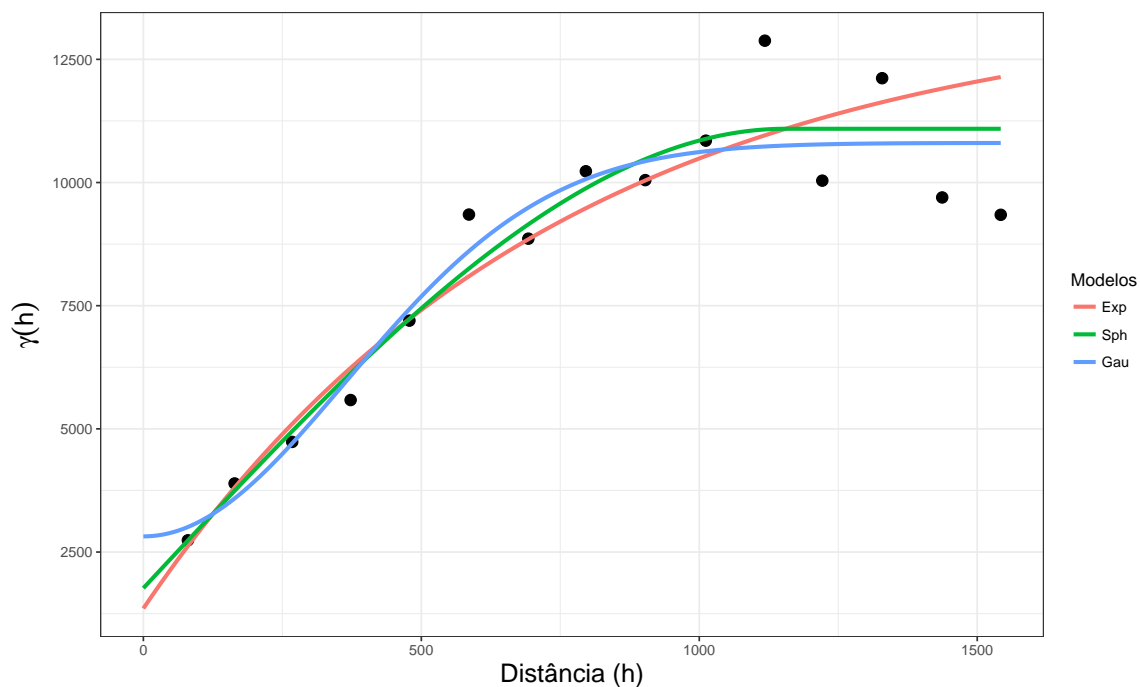


Figura 2: Semivariogramas ajustados

Foi considerado como critério de escolha a menor soma de quadrado dos erros. Observando a Tabela 1 temos que o modelo esférico apresentou menor soma de quadrado dos erros, sendo assim o modelo escolhido.

Tabela 1: Parâmetros dos modelos ajustados

Modelo	Exponencial	Esférico	Gaussiano
Alcance	737,87	1154,36	515,39
Efeito pepita	1354,61	1767,04	2817,48
Patamar	12307,76	9323,75	7985,83
Soma de quadrado	11287,385	7567,591	8147,89

Dessa forma, verificamos que a distância na qual o semivariograma atinge seu patamar é 1154 metros, ou seja, em distâncias superiores não há dependência espacial entre as observações.

Destacamos também que foi gerada uma malha de pontos a ser interpolada pela krigagem simples, krigagem ordinária, krigagem indicatriz e krigagem universal, ver Figura 3.

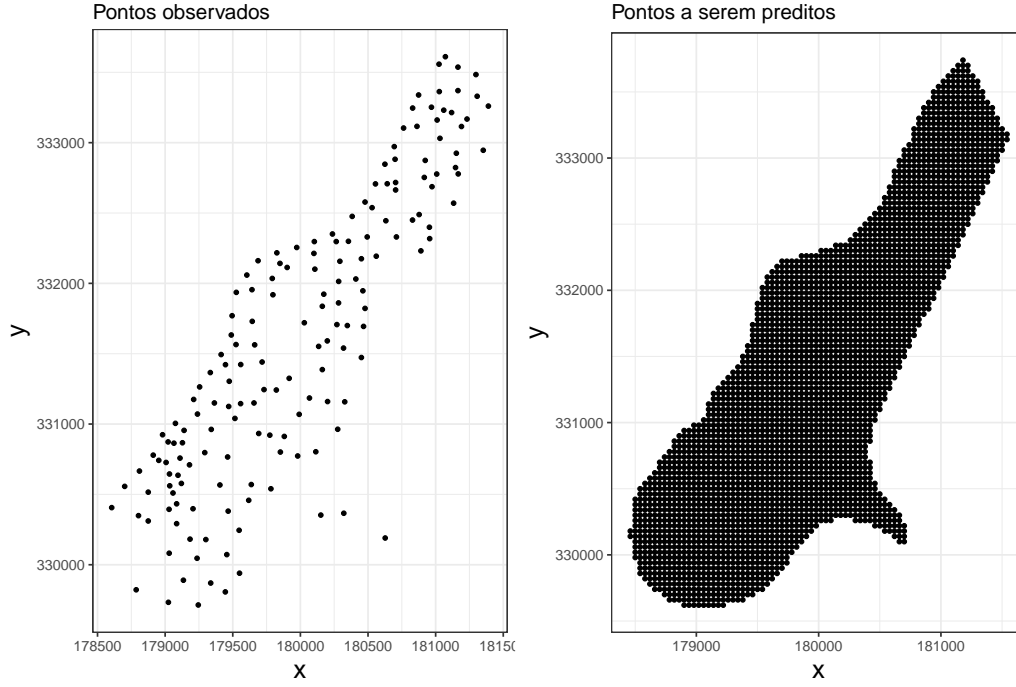


Figura 3: Pontos observados e malha de pontos gerada

3.2.2 Krigagem da Média

Uma medida amplamente utilizada para sumarizar uma variável aleatória é a média aritmética. Quando se calcula a média aritmética de alguma variável aleatória estamos estimando o parâmetro de locação μ da distribuição normal, isto é, estamos supondo que a variável aleatória tem distribuição normal com parâmetros μ e σ , sendo seus estimadores obtidos respectivamente por:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad (2)$$

e

$$s = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (3)$$

Por outro lado, quando a amostra possui localização no espaço assumimos que cada observação é uma variável aleatória, ou seja, temos um processo estocástico (1). Neste contexto, é interessante incorporar a dependência espacial entre as observações, no sentido de obter não mais um estimador e sim um preditor robusto para a média, em que cada observação é ponderada a partir da estrutura espacial do fenômeno.

Assumindo que

- (i) A média $\mu \in \mathbb{R}$ existe e é constante para todas localizações do domínio espacial D , isto é,

$$\mathbb{E}[Z(x)] = \mu \quad \forall x \in D,$$

- (ii) O processo estocástico é estacionário de segunda ordem com função de covariância dada por

$$C(h) = C[Z(x), Z(x+h)] = \mathbb{E}[Z(x)Z(x+h)] - \mu^2.$$

Matheron (1963) e Wackernagel (2013) definem um preditor linear m^* para a krigagem da média sendo uma combinação linear do processo estocástico avaliado em cada amostra x_i , isto é,

$$m^* = \sum_{i=1}^n w_i Z(x_i) \quad (4)$$

com ponderadores $w_i \in \mathbb{R}$ desconhecidos. Logo nosso objetivo é determinar os w_i , para isso iremos impor algumas restrições.

Não tendenciosidade

Inicialmente, queremos um preditor cujo erro $m^* - \mu$ seja em média 0, isto é, um preditor não tendencioso. Assim temos

$$\begin{aligned} \mathbb{E}[m^* - \mu] &= \mathbb{E}\left[\sum_{i=1}^n w_i Z(x_i) - \mu\right] = 0 \\ &= \sum_{i=1}^n w_i \mathbb{E}[Z(x_i)] - \mu = 0 \\ &= \mu \left[\sum_{i=1}^n w_i - 1\right] = 0 \\ \therefore \sum_{i=1}^n w_i &= 1. \end{aligned} \quad (5)$$

Portanto a soma dos ponderadores deve ser igual a um.

Variância mínima do erro de predição

Além disso, queremos que a variância do erro de predição seja mínima, uma vez que ela fornece informação sobre a acurácia do preditor linear m^* (LICHTENSTERN, 2013). Note que a variância do erro de predição é dada por

$$\begin{aligned} \mathbb{V}(m^* - \mu) &= \mathbb{E}\left[(m^* - \mu)^2\right] - \underbrace{(\mathbb{E}[m^* - \mu])^2}_{= \text{Vício}=0} \\ &= \mathbb{E}\left[(m^*)^2\right] - 2\mu \underbrace{\mathbb{E}[m^*]}_{=\mu} + \mu^2 \\ &= \mathbb{E}\left[\left(\sum_{i=1}^n w_i Z(x_i)\right)^2\right] - \mu^2 \\ &= \sum_{i=1}^n \sum_{j=1}^n w_i w_j \mathbb{E}[Z(x_i) Z(x_j)] - \mu^2 \sum_{i=1}^n w_i \sum_{j=1}^n w_j \\ &= \sum_{i=1}^n \sum_{j=1}^n w_i w_j \left(\mathbb{E}[Z(x_i) Z(x_j)] - \mu^2\right), \end{aligned}$$

Analisando a última expressão concluímos que

$$\mathbb{V}(m^* - \mu) = \sum_{i=1}^n \sum_{j=1}^n w_i w_j C(x_i - x_j). \quad (6)$$

Devemos determinar os ponderadores “ótimos” w_i que minimizam a Equação 6, sob a condição de não tendenciosidade Equação 5. Observe que temos um problema clássico de multiplicadores de Lagrange. Vamos considerar

$$f(\mathbf{w}) = \sum_{i=1}^n \sum_{j=1}^n w_i w_j C(\mathbf{x}_i - \mathbf{x}_j) \text{ e } g(\mathbf{w}) = \sum_{i=1}^n w_i = 1.$$

Assim, o método dos multiplicadores de Lagrange consiste em resolver o sistema de equações

$$\begin{cases} \nabla f(\mathbf{w}) = \lambda \nabla g(\mathbf{w}) \\ g(\mathbf{w}) = 1 \end{cases} \quad (7)$$

em que $\nabla f(\cdot)$ e $\nabla g(\cdot)$ são os gradientes da função $f(\cdot)$ e $g(\cdot)$, respectivamente.

Note que,

$$\begin{aligned} \nabla f(\mathbf{w}) &= \frac{\partial}{\partial w_i} f(w_i) = \frac{\partial}{\partial w_i} \sum_{i=1}^n w_i [w_1 C(\mathbf{x}_i - \mathbf{x}_1) + \dots + w_k C(\mathbf{x}_i - \mathbf{x}_k) + \dots + w_n C(\mathbf{x}_i - \mathbf{x}_n)] \\ &= w_1 C(\mathbf{x}_i - \mathbf{x}_1) + \dots + w_k C(\mathbf{x}_i - \mathbf{x}_k) + \dots + w_n C(\mathbf{x}_i - \mathbf{x}_n) \\ &= \sum_{j=1}^n w_j C(\mathbf{x}_i - \mathbf{x}_j), \quad i = 1, 2, \dots, n \end{aligned}$$

e

$$\lambda \nabla g(\vec{\omega}) = \lambda \frac{\partial}{\partial \omega_i} \sum_{i=1}^n \omega_i = \lambda \quad (8)$$

Dessa forma, o sistema linear com $n + 1$ equações pode ser expresso como

$$\begin{cases} \sum_{j=1}^n w_j C(\mathbf{x}_i - \mathbf{x}_j) - \lambda = 0 & \text{para } i = 1, 2, \dots, n \\ \sum_{j=1}^n w_j = 1 \end{cases} \quad (9)$$

Em forma matricial temos

$$\begin{bmatrix} C(\mathbf{x}_1 - \mathbf{x}_1) & C(\mathbf{x}_1 - \mathbf{x}_2) & \dots & C(\mathbf{x}_1 - \mathbf{x}_n) & -1 \\ C(\mathbf{x}_2 - \mathbf{x}_1) & C(\mathbf{x}_2 - \mathbf{x}_2) & \dots & C(\mathbf{x}_2 - \mathbf{x}_n) & -1 \\ & & \ddots & & \\ C(\mathbf{x}_n - \mathbf{x}_1) & C(\mathbf{x}_n - \mathbf{x}_2) & \dots & C(\mathbf{x}_n - \mathbf{x}_n) & -1 \\ 1 & 1 & \dots & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ \lambda \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \quad (10)$$

A solução desse sistema fornece os ponderadores ótimos w_i para prever o valor esperado do processo espacial $Z(\mathbf{x}_i)$, conforme Equação 4.

3.2.3 Krigagem Simples

Em diversas áreas do conhecimento o interesse da análise não se restringe na predição do valor esperado do processo espacial $Z(\mathbf{x})$, desejando também prever o valor de $Z(\mathbf{x})$ em algum ponto não amostrado \mathbf{x}_0 .

No caso da krigagem simples as suposições são (WACKERNAGEL, 2013):

- (i) A média $\mu \in \mathbb{R}$ é conhecida e constante para todas localizações do domínio espacial D , isto é,

$$\mathbb{E}[Z(x)] = \mu \quad \forall x \in D,$$

- (ii) O processo estocástico é estacionário de segunda ordem com função de covariância conhecida e dada por

$$C(h) = C[Z(x), Z(x+h)] = \mathbb{E}[Z(x) Z(x+h)] - \mu^2.$$

O preditor da krigagem simples utiliza a informação de cada amostra, o conhecimento do valor médio e o modelo teórico da covariância para prever algum ponto.

Conforme apresentado em ??) o preditor da krigagem simples $\hat{Z}(x_0)$ para a localização x_0 é dado por:

$$\hat{Z}(x_0) = \mu + \sum_{i=1}^n w_i [Z(x_i) - \mu] \quad (11)$$

com $w_i \in \mathbb{R}$ sendo os pesos desconhecidos de cada observação.

Agora iremos impor algumas restrições semelhantes a da krigagem da média para determinar os pesos w_i e então prever um ponto não amostrado.

Não tendenciosidade

Um preditor $\hat{Z}(x_0)$ de $Z(x_0)$ é dito ser não tendencioso se, e somente se, $\mathbb{E}[\hat{Z}(x_0) - Z(x_0)] = 0$. Observe que,

$$\begin{aligned} \mathbb{E}[\hat{Z}(x_0) - Z(x_0)] &= \mu + \sum_{i=1}^n w_i \underbrace{\mathbb{E}[Z(x_i) - \mu]}_{=0} - \underbrace{\mathbb{E}[Z(x_0)]}_{=\mu} \\ &= \mu - \mu = 0. \end{aligned}$$

Assim, não há restrição sobre os ponderadores $w_i, i = 1, 2, \dots, n$ imposta por essa propriedade.

Variância mínima do erro de predição

Além disso, em virtude da não tendenciosidade de $\hat{Z}(\mathbf{x}_0)$ a variância de predição é dada pelo erro quadrático médio, isto é:

$$\begin{aligned}
\mathbb{V}(\hat{Z}(\mathbf{x}_0) - Z(\mathbf{x}_0)) &= \mathbb{E} \left[(\hat{Z}(\mathbf{x}_0) - Z(\mathbf{x}_0))^2 \right] \\
&= \mathbb{E} \left[\left(\sum_{i=1}^n w_i (Z(\mathbf{x}_i) - \mu) + (\mu - Z(\mathbf{x}_0)) \right)^2 \right] \\
&= \sum_{i=1}^n \sum_{j=1}^n w_i w_j \mathbb{E} [(Z(\mathbf{x}_i) - \mu)(Z(\mathbf{x}_j) - \mu)] - 2 \sum_{i=1}^n w_i \mathbb{E} [(Z(\mathbf{x}_0) - \mu)(Z(\mathbf{x}_i) - \mu)] \\
&\quad + \mathbb{E} [(Z(\mathbf{x}_0) - \mu)^2] \\
&= \sum_{i=1}^n \sum_{j=1}^n w_i w_j \left(\mathbb{E}[Z(\mathbf{x}_i)Z(\mathbf{x}_j)] - \mu^2 \right) - 2 \sum_{i=1}^n w_i \left(\mathbb{E}[Z(\mathbf{x}_i)Z(\mathbf{x}_0)] - \mu^2 \right) \\
&\quad + \mathbb{V}(Z(\mathbf{x}_0)) \\
&= \sum_{i=1}^n \sum_{j=1}^n w_i w_j C(Z(\mathbf{x}_i), Z(\mathbf{x}_j)) - 2 \sum_{i=1}^n w_i C(Z(\mathbf{x}_i), Z(\mathbf{x}_0)) + C(Z(\mathbf{x}_0), Z(\mathbf{x}_0)) \\
&= C(\mathbf{0}) + \sum_{i=1}^n \sum_{j=1}^n w_i w_j C(\mathbf{x}_i, \mathbf{x}_j) - 2 \sum_{i=1}^n w_i C(\mathbf{x}_i, \mathbf{x}_0)
\end{aligned}$$

Podemos escrever a ultima expressão em termos matriciais como segue

$$\mathbb{V}(\hat{Z}(\mathbf{x}_0) - Z(\mathbf{x}_0)) = C(\mathbf{0}) + \mathbf{w}^T \Sigma \mathbf{w} - 2 \mathbf{w}^T \mathbf{c}_0 \quad (12)$$

Derivando e igualando a zero a Equação 12 obtemos:

$$\begin{aligned}
\frac{\partial}{\partial \mathbf{w}} \left(C(\mathbf{0}) + \mathbf{w}^T \Sigma \mathbf{w} - 2 \mathbf{w}^T \mathbf{c}_0 \right) &= (\Sigma + \Sigma^T) \mathbf{w} - 2 \mathbf{c}_0 = 0 \\
&= 2 \Sigma \mathbf{w} - 2 \mathbf{c}_0 = 0 \\
\therefore \Sigma \mathbf{w} &= \mathbf{c}_0.
\end{aligned}$$

Assim sendo, o sistema de equações para a krigagem simples, é dado por

$$\begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix} = \begin{pmatrix} C(x_1, x_1) & C(x_1, x_2) & \cdots & C(x_1, x_n) \\ C(x_2, x_1) & C(x_2, x_2) & \cdots & C(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ C(x_n, x_1) & C(x_n, x_2) & \cdots & C(x_n, x_n) \end{pmatrix}^{-1} \begin{pmatrix} C(x_1, x_0) \\ C(x_2, x_0) \\ \vdots \\ C(x_n, x_0) \end{pmatrix}$$

Após resolução deste sistema, obtemos os ponderadores \mathbf{w} e podemos interpolar qualquer realização do processo espacial $Z(x_i)$, conforme Equação 11.

3.2.4 Krigagem Ordinária

Em muitas situações a média do processo espacial é desconhecida, logo é importante estimar ela utilizando o modelo de semivariograma ajustado. Nesse sentido temos que a krigagem ordinária oferece uma solução otimizada quando a média do processo é desconhecida. Por esta razão a krigagem ordinária é o método mais utilizado na prática, quando o interesse é predizer um valor específico.

De acordo com Lichtenstern (2013), a krigagem ordinária segue as seguintes suposições:

- (i) A média $\mu \in \mathbb{R}$ é desconhecida para todas localizações do domínio espacial D , isto é,
- (ii) Os dados são provenientes de um processo estocástico intrinsecamente estacionário com função semivariância dada por:

$$\gamma(h) = \frac{1}{2} \mathbb{V}(Z(x+h) - Z(x)) = \frac{1}{2} \mathbb{E}[(Z(x+h) - Z(x))^2].$$

Já o preditor da krigagem ordinária, conforme apresentado em Diggle, Tawn e Moyeed (1998) $\hat{Z}(\mathbf{x}_0)$ para a localização \mathbf{x}_0 pode ser definido :

$$\hat{Z}(\mathbf{x}_0) = \sum_{i=1}^n w_i Z(\mathbf{x}_i). \quad (13)$$

com $w_i \in \mathbb{R}$ sendo os pesos desconhecidos de cada observação, isto é, fornecem a influencia de cada variável $Z(\mathbf{x}_i)$.

Vamos impor algumas restrições importante para encontrarmos os ponderados $\mathbf{w} = (w_1, w_2, \dots, w_n)$ e posteriormente predizermos algum ponto x_0 .

Não tendenciosidade

Queremos um preditor não viciado, isto é, que em média seja igual ao verdadeiro valor da variável aleatória. Assim temos:

$$\begin{aligned} \mathbb{E}[\hat{Z}(\mathbf{x}_0) - Z(\mathbf{x}_0)] &= \mathbb{E}\left[\sum_{i=1}^n w_i Z(\mathbf{x}_i) - Z(\mathbf{x}_0) \sum_{i=1}^n w_i\right] \\ &= \sum_{i=1}^n w_i \mathbb{E}[Z(\mathbf{x}_i) - Z(\mathbf{x}_0)] = 0 \\ \therefore \sum_{i=1}^n w_i &= 1. \end{aligned} \quad (14)$$

Note, que a soma dos ponderadores deve ser igual a um.

Variância mínima do erro de predição

Outra importante propriedade é a variância do erro de predição, queremos que ela seja a menor possível. Vamos primeiramente encontra-la, Pelo fato de \hat{Z}_0 ser não tendencioso, então a variância do erro de estimação é igual ao erro quadrático médio, assim temos:

$$\begin{aligned} \mathbb{V}(\hat{Z}(\mathbf{x}_0) - Z(\mathbf{x}_0)) &= \mathbb{E}[(\hat{Z}(\mathbf{x}_0) - Z(\mathbf{x}_0))^2] \\ &= \mathbb{E}\left[\left(\sum_{i=1}^n w_i Z(\mathbf{x}_i) - Z(\mathbf{x}_0) \sum_{i=1}^n w_i\right)^2\right] \\ &= \sum_{i=1}^n \sum_{j=1}^n w_i w_j \mathbb{E}[Z(\mathbf{x}_i)Z(\mathbf{x}_j)] - 2 \sum_{i=1}^n w_i \mathbb{E}[Z(\mathbf{x}_i)Z(\mathbf{x}_0)] + \mathbb{E}[(Z(\mathbf{x}_0))^2] \\ &= - \sum_{i=1}^n \sum_{j=1}^n w_i w_j \frac{\mathbb{E}[(Z(\mathbf{x}_i) - Z(\mathbf{x}_j))^2]}{2} + 2 \sum_{i=1}^n w_i \frac{\mathbb{E}[(Z(\mathbf{x}_i) - Z(\mathbf{x}_0))^2]}{2}. \end{aligned}$$

Utilizando a definição da função semivariância, obtemos a variância de predição dada por:

$$\mathbb{V}(\hat{Z}(x_0) - Z(x_0)) = - \sum_{i=1}^n \sum_{j=1}^n w_i w_j \gamma(x_i - x_j) + 2 \sum_{i=1}^n w_i \gamma(x_i - x_0). \quad (15)$$

Agora devemos minimizar (15) sob a condição de não tendenciosidade (14) e assim encontrar os ponderados ótimos \mathbf{w} . Observe que temos um problema clássico de multiplicadores de Lagrange. Dessa forma, minimizando a Equação 15, sujeita a restrição da Equação 14, obtemos um sistema linear com $n + 1$ equações e $n + 1$ incógnitas:

$$\underbrace{\begin{bmatrix} \gamma(x_1 - x_1) & \gamma(x_1 - x_2) & \cdots & \gamma(x_1 - x_n) & 1 \\ \gamma(x_2 - x_1) & \gamma(x_2 - x_2) & \cdots & \gamma(x_2 - x_n) & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \gamma(x_n - x_1) & \gamma(x_n - x_2) & \cdots & \gamma(x_n - x_n) & 1 \\ 1 & 1 & \cdots & 1 & 0 \end{bmatrix}}_{\Sigma} \cdot \underbrace{\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ \lambda \end{bmatrix}}_{\Omega} = \underbrace{\begin{bmatrix} \gamma(x_0 - x_1) \\ \gamma(x_0 - x_2) \\ \vdots \\ \gamma(x_0 - x_n) \\ 1 \end{bmatrix}}_{\Gamma}$$

Assim temos que:

$$\Sigma \cdot \Omega = \Gamma \rightarrow \Omega = \Sigma^{-1} \cdot \Gamma \quad (16)$$

Note que para obtenção dos ponderadores de krigagem é necessário que a matriz Σ seja positiva definida, pois precisamos de sua inversa. Finalmente, após determinação dos ponderados \mathbf{w} podemos interpolar qualquer realização do processo espacial, conforme a Equação 4, que é uma estimativa do valor esperado da variável no ponto x_0 , isto é, $\mathbb{E}[Z(x_0)]$

3.2.5 Krigagem Indicatriz

Como a krigagem ordinária fornece apenas estimativa da esperança da variável aleatória $Z(x_i)$ em pontos desconhecidos, as possibilidades de interpretação e uso destes resultados são limitadas. Uma das alternativas, que possibilita a estimação não só da esperança, mas de toda a função de distribuição acumulada da variável em cada ponto é a krigagem indicatriz (ALMEIDA; JR, 1996).

Em sua essência a krigagem indicatriz se baseia na transformação do conjunto de dados em variáveis indicadoras, mantendo-se a metodologia e os pressupostos da krigagem ordinária. Consideraremos a transformação indicadora para um dado ponto de corte z_c .

$$I(x_i | z_c) = \begin{cases} 1, & \text{se } Z(x_i) \leq z_c \\ 0, & \text{caso contrário.} \end{cases} \quad (17)$$

Conforme é apresentado em Isaacs e Srivastava (1989) temos que a esperança e a variância de cada variável indicadora $I(x_i | z_c)$ podem ser obtidas respectivamente por:

$$\mathbb{E}[I(x_i | z_c)] = 1 \times P[Z(x_i) \leq z_c] + 0 \times P[Z(x_i) > z_c] = P[Z(x_i) \leq z_c] \quad (18)$$

e

$$\mathbb{V}[I(x_i | z_c)] = P[Z(x_i) \leq z_c] \times \{1 - P[Z(x_i) \leq z_c]\} \quad (19)$$

Como pode-se observar a estima um ponto da função de distribuição acumulada condicional de $Z(x_i)$ em z_c . Além disso, conforme define (JOURNEL, 1980) o preditor da krigagem indicadora $\hat{I}(x_0 | z_c)$ para a localização x_0 é dado por:

$$\hat{I}(x_0 | z_c) = \sum_{i=1}^n w_i I(x_i | z_c) \quad (20)$$

em que $\mathbf{w} = w_1, \dots, w_n$ os ponderadores devem ser estimados de modo a garantir as propriedades de variância mínima e não tendenciosidade, tal como na krigagem ordinária (ver Equações (14) e (15)).

Assim sendo repetindo os mesmos procedimentos da krigagem ordinária resultamos em um sistema com $n + 1$ equações e $n + 1$ incógnitas:

$$\underbrace{\begin{bmatrix} \gamma(x_1 - x_1) & \gamma(x_1 - x_2) & \cdots & \gamma(x_1 - x_n) & 1 \\ \gamma(x_2 - x_1) & \gamma(x_2 - x_2) & \cdots & \gamma(x_2 - x_n) & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \gamma(x_n - x_1) & \gamma(x_n - x_2) & \cdots & \gamma(x_n - x_n) & 1 \\ 1 & 1 & \cdots & 1 & 0 \end{bmatrix}}_{\Sigma} \cdot \underbrace{\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ \lambda \end{bmatrix}}_{\Omega} = \underbrace{\begin{bmatrix} \gamma(x_0 - x_1) \\ \gamma(x_0 - x_2) \\ \vdots \\ \gamma(x_0 - x_n) \\ 1 \end{bmatrix}}_{\Gamma}$$

Logo têm-se:

$$\Sigma \cdot \Omega = \Gamma \rightarrow \Omega = \Sigma^{-1} \cdot \Gamma \quad (21)$$

Ou seja, após a determinação da matriz Ω podemos prever a função acumulada de um ponto condicionada a um dado ponto de corte z_c . É importante ressaltar que o semivariograma experimental para a krigagem indicatriz são obtidos por (ISAACS; SRIVASTAVA, 1989):

$$\gamma(h | z_c) = \frac{1}{2N(h)} \sum_{i=1}^{N(h)} [I(x_i | z_c) - I(x_i + h | z_c)]^2 \quad (22)$$

Sendo:

- h a distância entre as observações;
- $N(h)$ o número de pares de valores medidos, separados por uma distância h ;
- z_c o ponto de corte.

Portanto, a única diferença entre a krigagem ordinária e a krigagem indicatriz é a natureza da variável aleatória, sendo que na krigagem ordinária predizemos valores de uma variável aleatória contínua, ao passo que na krigagem indicatriz predizemos valores de uma variável aleatória dicotômica.

Por fim, como apontado por Almeida e Jr (1996) ao se estimar a função de distribuição acumulada podemos obter mais informações sobre a área de estudo, sendo possível a construção de mapas não só de médias, mas de outras medidas tais como: mediana, quantis e probabilidades.

3.2.6 Krigagem Universal

Além da média do processo espacial ser desconhecida, em diversas aplicações não é razoável supor que a média independe da localização espacial. Particularmente

isto ocorre quando a análise abrange grandes regiões. O método de predição espacial adequado nestes casos é a krigagem universal.

Conforme discutido por Santos (2010) o preditor da Krigagem universal é uma proposta para situações onde a função média do processo é desconhecida e ainda não estacionária. De acordo com Lichtenstern (2013), a krigagem universal segue as seguintes suposições:

1. Assume-se que $Z(\mathbf{x})$ é decomposta em um termo função determinística $\mu(\mathbf{x})$ e uma função aleatória $Y(\mathbf{x})$, tal que:

$$Z(\mathbf{x}) = \mu(\mathbf{x}) + Y(\mathbf{x})$$

2. Supõe-se que $Y(\mathbf{x})$ seja intrinsecamente estacionária com média zero e função variograma $\gamma_Y(\mathbf{h})$, denominada função variograma residual de $Z(\mathbf{x})$, isto é, $\forall \mathbf{x}, \mathbf{x} + \mathbf{h} \in D$:

$$\begin{aligned} \mathbb{E}[Z(\mathbf{x})] &= \mu(\mathbf{x}) \\ \gamma_Y(\mathbf{h}) &= \frac{1}{2} \mathbb{V}[Y(\mathbf{x} + \mathbf{h}) - Y(\mathbf{x})] = \frac{1}{2} \mathbb{E}[(Y(\mathbf{x} + \mathbf{h}) - Y(\mathbf{x}))^2] \end{aligned} \quad (23)$$

3. Por fim, seja f_0, f_1, \dots, f_L as funções determinísticas de uma coordenada geográfica $\mathbf{x} \in D$ com $L \in \mathbb{N}$ um número conhecido e selecionáveis de funções básicas $f_l : D \rightarrow \mathbb{R}$, $l = 0, \dots, L$. Assumimos que $\mu(\mathbf{x})$ é uma combinação linear dessa funções avaliadas em \mathbf{x}

$$\mu(\mathbf{x}) = \sum_{l=0}^L a_l f_l(\mathbf{x}) \quad (24)$$

com coeficientes desconhecidos $a_l \in \mathbb{R}$.

O preditor da krigagem universal $\hat{Z}(\mathbf{x}_0)$ para o valor de $Z(\mathbf{x})$ a um certo ponto \mathbf{x}_0 é a soma linear

$$\hat{Z}(\mathbf{x}_0) = \sum_{i=1}^n w_i Z(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{Z} \quad (25)$$

com pesos $w_i \in \mathbb{R}$, $i = 1, \dots, n$ correspondente a cada função aleatória $Z(\mathbf{x})$ no ponto amostral \mathbf{x}_i e $\mathbf{w} = (w_1, \dots, w_n)^\top \in \mathbb{R}^n$.

As condições de não tendenciosidade do preditor bem como variância mínima do erro de predição são complexas de serem demonstradas e não serão apresentadas neste trabalho. No entanto leitores interessados podem encontrar as demonstrações nos livros de Isaacs e Srivastava (1989), Diggle, Tawn e Moyeed (1998) ou nas recentes dissertações de Santos (2010) e Lichtenstern (2013).

3.3 Resultados e Discussões

Nas próximas subseções, serão apresentados e discutidos os resultados da aplicação de cada krigagem utilizando o conjunto de dados **meuse**.

3.3.1 Krigagem da Média no R

Como aplicação predizemos o valor esperado, através de uma função implementada junto ao software R. No código a seguir apresentamos a função e o comando comparando a krigagem da média com a média aritmética.

Código – R 3.1: Função para a krigagem da média

```
bib <- c("gstat", "sp", "geoR")
sapply(bib, require, character.only = TRUE)

krigagem.media <- function(data, coords, model, pars,
  nug = 0){
  nr <- nrow(coords)
  cov.mat <- varcov.spatial(coords=coords,
    cov.model=model, cov.pars=pars, nugget = nug)$varcov
  c <- c(rep(-1,l=nr),0)
  l <- c(rep(1,l=nr))
  cov.mat <- rbind(cov.mat,l)
  cov.mat <- cbind(cov.mat,c)
  R <- matrix(c(rep(0,l=nr),1))
  wi <- (solve(cov.mat)%*%R)[-nr]
  mean.krige <- sum(data*wi)
  return(mean.krige)
}

data("meuse.all")
coordinates(meuse.all) <- ~x+y
vemp <- variogram(lead ~ 1, data=meuse.all, cressie = T)
vfit <- fit.variogram(vemp, model = vgm("Sph"),
  fit.method = 1)

krigagem.media(data = meuse.all$lead, coords =
  meuse.all@coords, model = "sph",
  pars = c(vfit$psill[2], vfit$range[2]),
  nug = vfit$psill[1])
mean(meuse.all$lead)
```

É possível observar pela Tabela 2 uma grande distinção entre o valor médio estimado supondo que os dados provêm de uma distribuição normal e assumindo que o cobre é um processo estocástico com estrutura de dependência espacial descrita pelo modelo de semivariograma ajustado. Sendo assim, podemos inferir que o valor médio de cobre encontrado nesta localização é de 180,77 mg/kg.

Tabela 2: Comparação entre a média aritmética e krigagem da média

Média aritmética	Krigagem da média
148,5549	180,7796

3.3.2 Krigagem Simples no R

Tendo em vista o modelo de semivariograma ajustado e descrito na seção anterior, empregamos a krigagem simples para prever a malha de pontos da Figura 3. Pode-se observar pela Figura 4 a concentração de chumbo visa ser maior nas margens do rio Meuse, em contrapartida verifica-se uma baixa concentração no centro da malha de pontos.

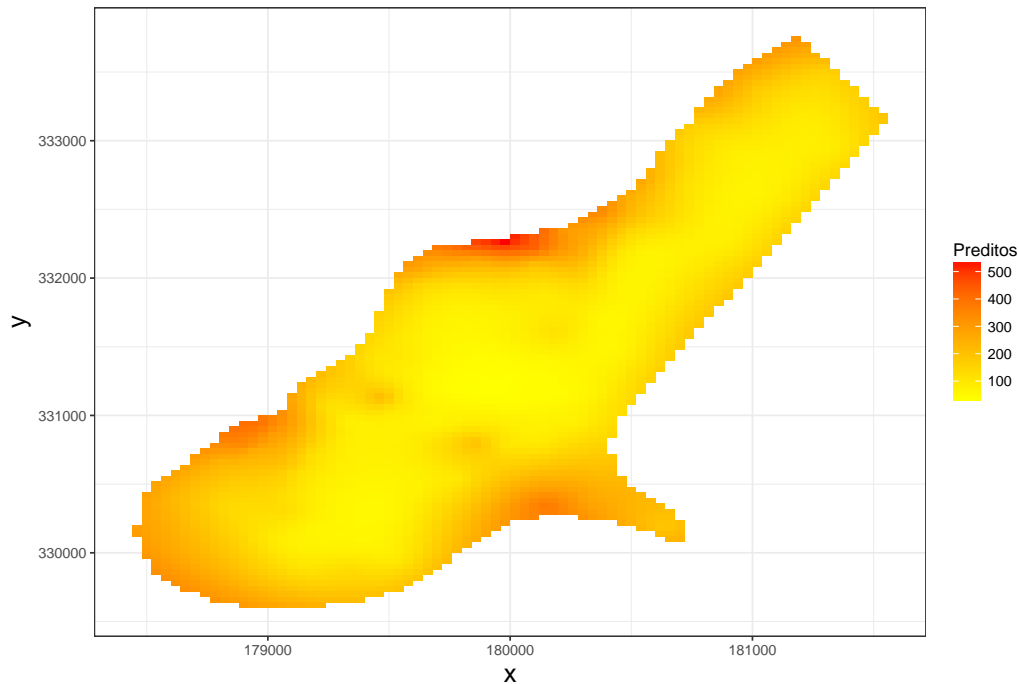


Figura 4: Malha de pontos interpolados pela krigagem simples

É apresentado na sequência o comando para realização da krigagem simples utilizando o pacote `gstat` (PEBESMA, 2007).

Código – R 3.2: Comandos para krigagem simples

```
library(gstat)

data("meuse.all")
data("meuse.all")
coordinates(meuse.all) <- ~ x + y
coordinates(meuse.grid) <- ~ x + y

mu <- mean(meuse.all@data$zinc)
preditos_kgs <- krige(lead ~ 1, meuse.all, meuse.grid,
  model=vfit, beta=mu)
head(preditos_kgs)
```

3.3.3 Krigagem Ordinária no R

Considerando o modelo de semivariograma ajustado e descrito nas seções anteriores, aplicamos a krigagem ordinária com intuito de prever a malha de pontos

da Figura 3. Como pode ser visto pela Figura 5 não há grandes diferenças, neste caso, entre a krigagem simples e ordinária, sendo que a concentração de chumbo é maior nas margens do rio Meuse, e é verificado uma baixa concentração no centro.

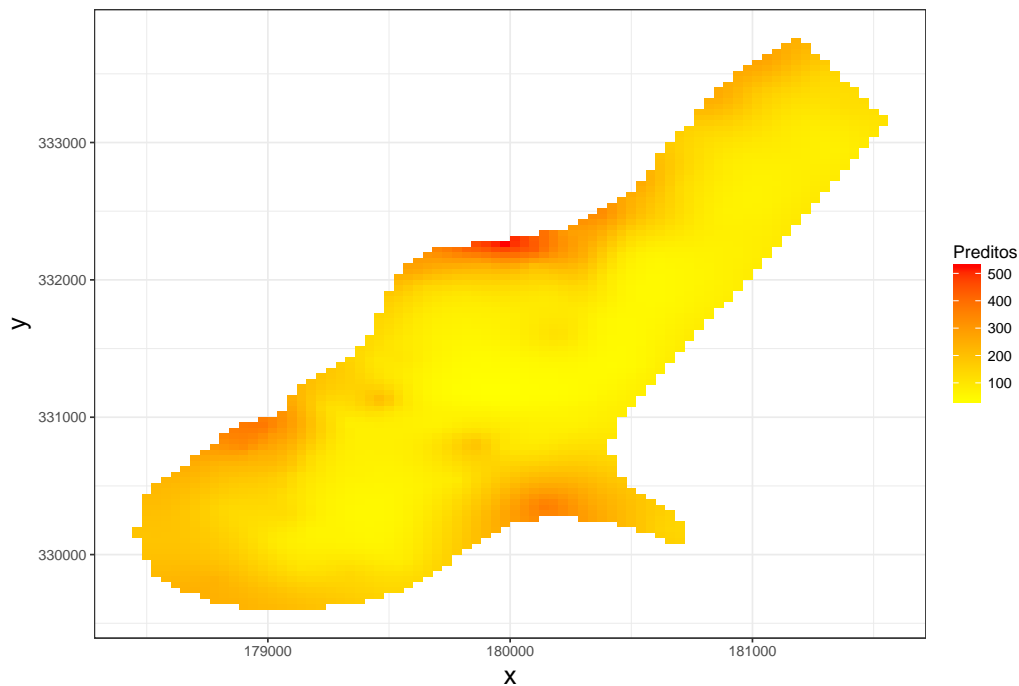


Figura 5: Malha de pontos interpolados pela krigagem ordinária

Em seguida, apresentamos os comandos para realização da krigagem ordinária utilizando a biblioteca `gstat` (PEBESMA, 2007).

Código – R 3.3: Comandos para krigagem ordinária

```
library(gstat)

data("meuse.all")
data("meuse.grid")
coordinates(meuse.all) <- ~ x + y
coordinates(meuse.grid) <- ~ x + y

preditos_kgo <- krige(lead ~ 1, meuse.all, meuse.grid,
  model=vfit)
head(preditos_kgo)
```

3.3.4 Krigagem Indicatriz no R

De acordo com o Atlas Ambiental Digital de Berlim¹ o consumo de vegetais folhosos devem ser evitados quando o nível de chumbo (Pb) esta acima de 100 mg/kg. Nesse sentido adotaremos como ponto de corte $z_c = 100$, assim iremos

¹<http://www.stadtentwicklung.berlin.de/umwelt/umweltatlas/ed103103.htm>

mostrar como ajustar um modelo de semivariograma para a krigagem indicatriz e posteriormente prever a malha de pontos da Figura 3.

A partir da Figura 6 observamos que o modelo esférico ajustado com patamar 0,20434 metros, alcance de 687,124 metros e efeito pepita de 0,06507 metros compreende razoavelmente o semivariograma empírico.

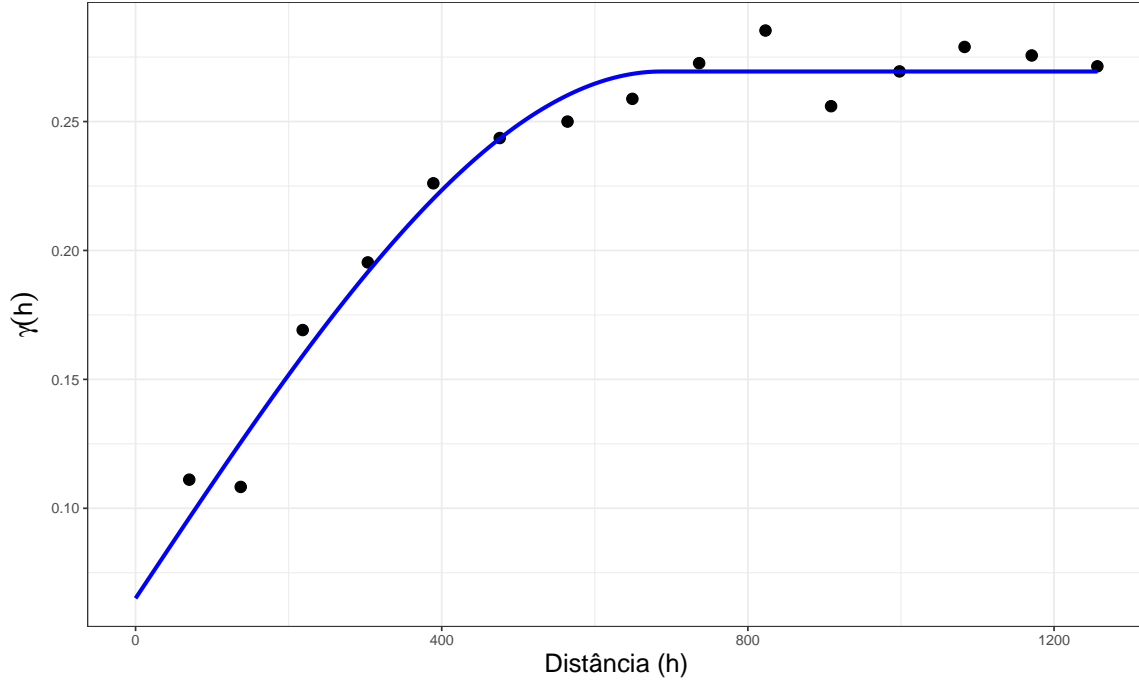


Figura 6: Semivariograma empírico e ajustado.

Tendo em vista o modelo de semivariograma ajustado, a krigagem indicatriz foi empregada com intuito de estimar as probabilidades de encontrar valores superiores a 100 mg/kg de chumbo na região da Figura 3.

Na Figura 7 temos um mapa de probabilidade, isto é, a probabilidade estimada de observarmos valores superiores a $z_c = 100$ mg/kg de chumbo. Como pode-se verificar é muito provável que nas margens do rio Meuse o nível de chumbo seja superior a 100 mg/kg. Em contrapartida, no centro da malha é esperado uma probabilidade baixa, próximo de zero, de haver níveis de chumbo acima de 100 mg/kg.

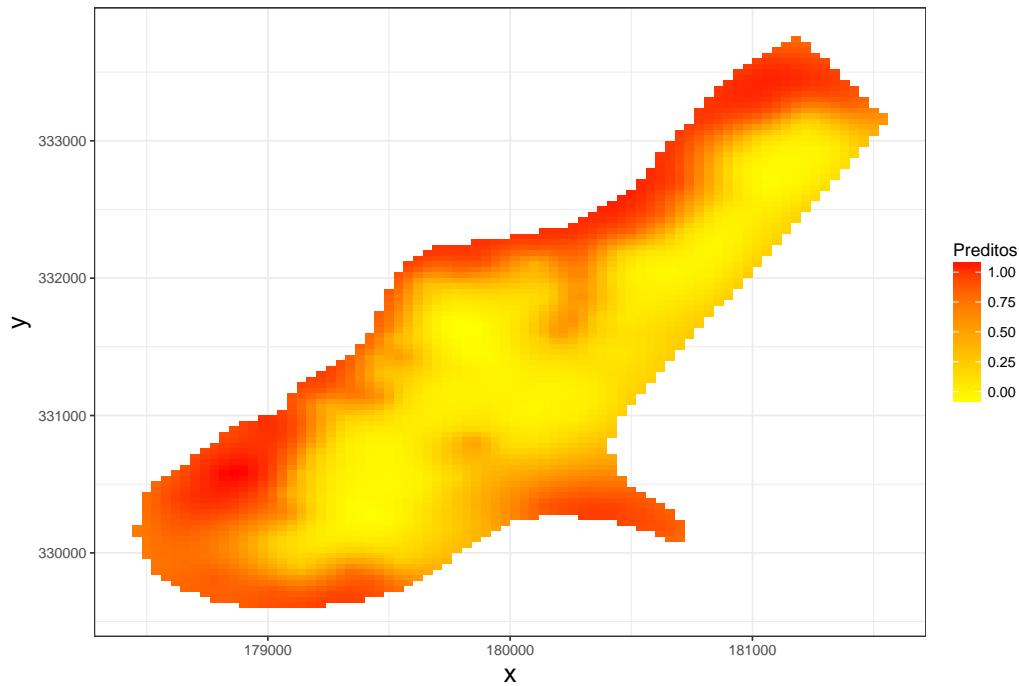


Figura 7: Malha de pontos interpolados pela krigagem indicatriz

Na sequência, apresentamos os comandos utilizados para realização da krigagem indicatriz, por meio da biblioteca `gstat` (PEBESMA, 2007).

Código – R 3.4: Comandos para krigagem indicatriz

```
library(gstat)
data("meuse.all"); data("meuse.grid")
coordinates(meuse.all) <- ~ x+y
coordinates(meuse.grid) <- ~ x+y
meuse.all$lead.i <- meuse.all$lead > 100

vi <- variogram(lead.i ~ 1, location = meuse.all,
  cutoff = 1300)
vimf <- fit.variogram(vi, vgm(0.12, "Sph", 1300, 0))

preditos_kgi <- krige(lead.i ~ 1, loc=meuse.all,
  newdata=meuse.grid, model=vimf)
head(preditos_kgi)
```

3.3.5 Krigagem Universal no R

Para fins de ilustração iremos considerar a média das concentrações (mg/kg) de chumbo do rio Meuse na Holanda, como uma função linear da latitude e longitude observadas. Para isso, teremos que ajustar outro semivariograma considerando a latitude e longitude no modelo e depois realizar as predições, com base na krigagem universal.

Na Figura 6 temos o semivariograma empírico e teórico, considerando o modelo adotado. Podemos notar que o modelo esférico ajustado com patamar

7062,392 metros, alcance de 1367,372 metros e efeito pepita de 1962.532 metros se ajusta bem ao semivariograma empírico.

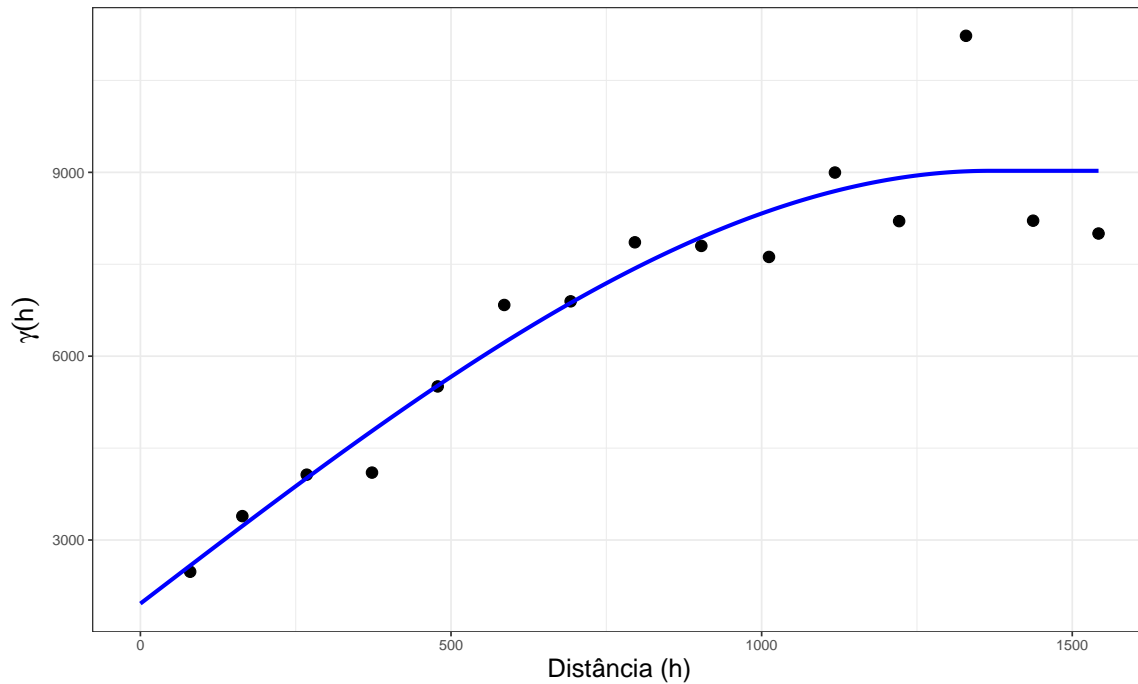


Figura 8: Semivariograma empírico e ajustado para a krigagem universal.

Tendo em vista o modelo de semivariograma ajustado, a krigagem universal foi utilizada para prever a malha de pontos exposta na Figura 3. A Figura 9 exibe a malha de pontos interpolados pela krigagem universal. Pode-se observar que a concentração de chumbo é maior nas margens do rio, ao passo que aparenta ter concentração menor de chumbo no centro do mapa.

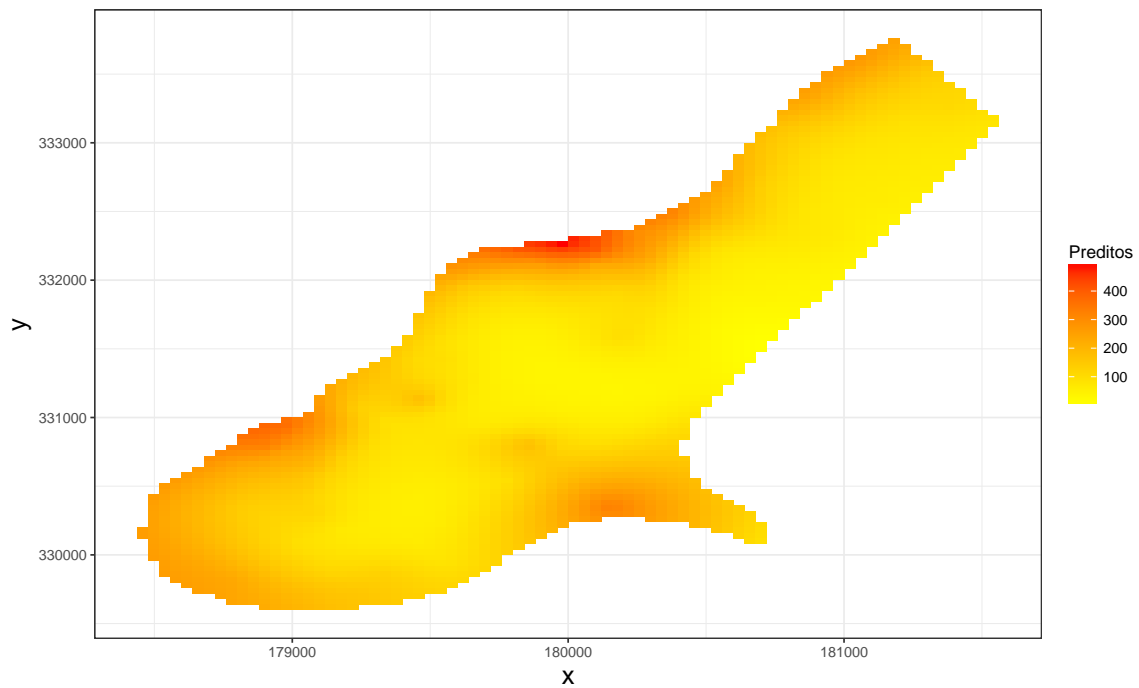


Figura 9: Malha de pontos interpolados pela krigagem universal.

Na sequência apresenta-se os comandos utilizados para realizar a krigagem universal utilizando o software R e a biblioteca **gstat** (PEBESMA, 2007).

Código – R 3.5: Comandos para krigagem universal

```
library(gstat)
data("meuse.all");data("meuse.grid")
coordinates(meuse.all) <- ~ x+y
coordinates(meuse.grid) <- ~ x+y

vi <- variogram(lead.i ~ x + y, location = meuse.all,
  cutoff = 1300)
vimf <- fit.variogram(vi, model = vgm(c("Exp", "Sph",
  "Gau"))))

preditos_kgu <- krige(lead ~ x + y, loc=meuse.all,
  newdata=meuse.grid, model=vimf)
head(preditos_kgu)
```

4 Conclusão

Ao término deste projeto de iniciação científica iremos ressaltar alguns aspectos importantes. Inicialmente, destacamos que a realização das propostas estabelecidas foi cumprida dentro do prazo definido. O segundo aspecto refere-se aos novos conhecimentos adquiridos pelo acadêmico, sendo que em sua maioria não são apresentados na graduação.

Por fim, no que tange aspectos relacionados ao tema estudado destacamos os seguintes: (i) as demonstrações presente neste relatório podem colaborar para melhor entendimento didático dos tipos de krigagem; (ii) uma função para realização da krigagem da média foi proposta, uma vez que foi o único método estudado que não conseguimos encontrar sua implementação no ambiente estatístico R; (iii) foi utilizado o pacote `gstat` para aplicação das demais krigagem, no entanto ressaltamos que os pacotes `geoR` (JR; DIGGLE, 2001) e `RandomFields` (SCHLATHER et al., 2016) também possuem funções para aplicação da krigagem; (iv) recomenda-se a utilização da krigagem ordinária para predição de valores desconhecidos quando a suposição de estacionaridade é garantida; (v) em contrapartida recomenda-se a utilização da krigagem universal quando a média do processo espacial varia por localização ou algum outro efeito aleatório presente no estudo.

Referências

- ALMEIDA, C. F. P.; JR, P. J. R. **Estimativa da distribuicao espacial de retencao de água em um solo utilizando krigagem indicatriz**. 1996.
- CRESSIE, N. Fitting variogram models by weighted least squares. **Journal of the International Association for Mathematical Geology**, Springer, v. 17, n. 5, p. 563–586, 1985.
- DIGGLE, P. J.; TAWN, J.; MOYEED, R. Model-based geostatistics. **Journal of the Royal Statistical Society: Series C (Applied Statistics)**, Wiley Online Library, v. 47, n. 3, p. 299–350, 1998.
- ISAACS, E. H.; SRIVASTAVA, R. M. **Applied geostatistics**. Oxford University Press, 1989.
- JOURNEL, A. The lognormal approach to predicting local distributions of selective mining unit grades. **Journal of the International Association for Mathematical Geology**, Springer, v. 12, n. 4, p. 285–303, 1980.
- JR, P. J. R.; DIGGLE, P. J. geor: a package for geostatistical analysis. **R news**, London, v. 1, n. 2, p. 14–18, 2001.
- LICHTENSTERN, A. **Kriging methods in spatial statistics**. Tese (Doutorado) — Technische Universitat Munchen, 2013.
- MATHERON, G. Principles of geostatistics. **Economic geology**, Society of Economic Geologists, v. 58, n. 8, p. 1246–1266, 1963.
- PEBESMA, E. J. The gstat package. **R package version 0.9-42**, 2007.
- SANTOS, G. R. dos. **Hierarquizacao Geometrica dos Preditores Geoestatisticos**. Tese (Doutorado) — Universidade Federal de Lavras - UFLA, 2010.
- SCHLATHER, M. et al. **RandomFields: Simulation and Analysis of Random Fields**. 2016. R package version 3.1.36.
- WACKERNAGEL, H. **Multivariate geostatistics: an introduction with applications**. Springer Science & Business Media, 2013.