

Estatística Espacial Geoestatística I

Raquel Menezes

Departamento de Matemática
Universidade do Minho

Setembro de 2023



Lembrar a importância da AED

- A análise exploratória de dados (AED) deve cobrir quer os aspetos **não espaciais** (e.g. boxplot) quer os **espaciais** (e.g. mapas).
Numa análise exploratória devemos fazer análise espacial e não espacial
- Atenção, as estatísticas pontuais não são suficientes!

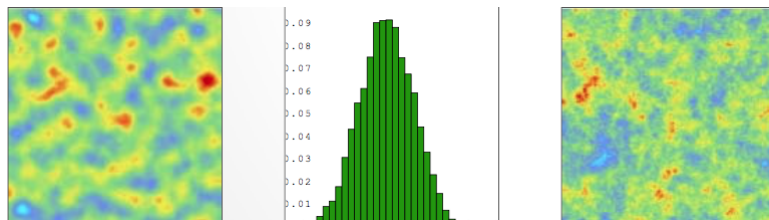


Figure: Dois Campos Aleatórios - CA distintos (painéis da esquerda e da direita), no entanto partilham o mesmo histograma (painel do centro).

Nota: Na terminologia inglesa um Campo Aleatório (CA) diz-se um *Random Field* (RF).

Dados referentes a pontos - Geoestatística

- O nosso interesse ir-se-á focar na
 - ▷ modelação de dados y_1, \dots, y_n recolhidos em diferentes localizações x_1, \dots, x_n (por exemplo, latitude e longitude) na região de estudo $D \subset \mathbb{R}^2$
 - ▷ **estimação da estrutura de correlação espacial** subjacente aos dados observados
- Iremos assumir um **domínio espacial contínuo**
 - ▷ processo estocástico $Y(x)$ pode ser medido em $\forall x \in D \subset \mathbb{R}^2$
 - ▷ classicamente referido por **processo geoestatístico** (razões históricas)
 - ▷ MAS atualmente cobrindo as diversas áreas de aplicação

Dados espaciais versus séries temporais

Semelhanças:

- Ambos consideram um conjunto discreto de observações, nos tempos t_1, \dots, t_n de $Y(t)$ ou nas localizações x_1, \dots, x_n de $Y(x)$, para estimar a estrutura de autocorrelação subjacente ao processo¹

Diferenças:

- Os dados nas **séries temporais têm apenas uma direção**, que é do passado para o tempo mais recente, podendo assim ser ordenados na reta do tempo

Aqui no espaço não temos uma ordem associada.

FAC



¹Nos dados espaciais, em alternativa à função autocorrelação, recorre-se frequentemente à **função variograma**

Notação

- $Y(\mathbf{x})$ identifica o fenômeno em estudo (por exemplo, abundância de uma espécie ou chuva) que **depende de coordenadas espaciais** $\mathbf{x} \in D \subset \mathbb{R}^2$, onde D é uma região limitada
- Se temos as localizações espaciais $\mathbf{x}_1, \dots, \mathbf{x}_n$, então $Y(\mathbf{x}_1), \dots, Y(\mathbf{x}_n)$ identifica dados observados nessas localizações. Essas observações podem ser obtidas a partir de uma ou mais variáveis discretas ou contínuas.

- $Y(\mathbf{x})$ é geralmente definido através da distribuição de dimensão finita

$$F_{\mathbf{x}_1, \dots, \mathbf{x}_n}(y_1, \dots, y_n) = P\{Y(\mathbf{x}_1) \leq y_1, \dots, Y(\mathbf{x}_n) \leq y_n\}, \quad n \geq 1$$

5 / 26

Estacionariedade

- $Y(\mathbf{x})$ é **estrita ou fortemente estacionário** se, para qualquer vector $\mathbf{u} \in \mathbb{R}^2$, temos

$$F_{\mathbf{x}_1, \dots, \mathbf{x}_n}(y_1, \dots, y_n) = F_{\mathbf{x}_1 + \mathbf{u}, \dots, \mathbf{x}_n + \mathbf{u}}(y_1, \dots, y_n), \quad \forall n, \mathbf{u}$$

o que significa que $Y(\mathbf{x})$ **permanece invariante quando sujeito a transformações de translação de suas coordenadas**

- $Y(\mathbf{x})$ é **estacionário de segunda-ordem** (ou fracamente estacionário), se

$$E[Y(\mathbf{x})] = \mu(\mathbf{x}) = \mu \quad \forall \mathbf{x} \in D$$

$$\text{Cov}[Y(\mathbf{x}_i), Y(\mathbf{x}_j)] = c(\mathbf{x}_i - \mathbf{x}_j), \quad \forall \mathbf{x}_i, \mathbf{x}_j \in D$$

$c(\cdot)$ é chamada função de **covariância** estacionária e $\mu(\mathbf{x})$ é conhecido como **tendência** do processo

- $Y(\mathbf{x})$ é **intrinsecamente estacionário**, se ²

$$E[Y(\mathbf{x})] = \mu(\mathbf{x}) = \mu \iff E[Y(\mathbf{x}_i) - Y(\mathbf{x}_j)] = 0$$

$$\text{Var}[Y(\mathbf{x}_i) - Y(\mathbf{x}_j)] = 2\gamma(\mathbf{x}_i - \mathbf{x}_j), \quad \forall \mathbf{x}_i, \mathbf{x}_j \in D$$

$\gamma(\cdot)$ é chamado **variograma**

²Semelhante à série temporal, se Y_t não for estacionário, pode-se considerar as diferenças (de primeira ordem) entre os tempos

7 / 26

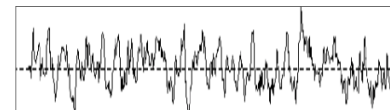
Hipóteses

- $Y(\mathbf{x})$ **estacionário**: **invariância sob translações/rotações no espaço**.

Se for feita uma translação (vector \mathbf{u}) a distribuição conjunta matém-se.

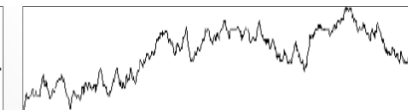
- $Y(\mathbf{x})$ **estacionário de ordem 2**: os dois primeiros momentos existem e são invariantes sob translações/rotações. Média, variância e covariância têm de ser constantes, não podem depender de \mathbf{x} . Os momentos (média, variância e covariância) são os mesmos quando fazemos uma translação.

- $Y(\mathbf{x})$ **intrínseco**: os incrementos são estacionários de ordem 2.



Estacionário

Parece ser estacionário na média e na variância



Intrínseco

Não parece ser estacionário na média

6 / 26

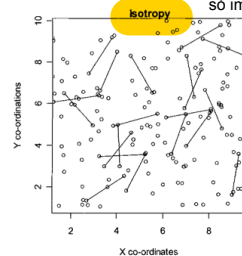
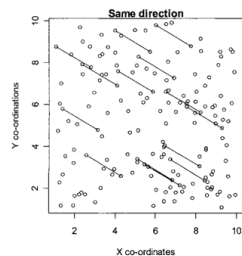
Hipóteses

- Um processo $Y(\mathbf{x})$ estacionário também é intrínseco, mas um $Y(\mathbf{x})$ intrínseco nem sempre é estacionário.
- O foco não está na média... mas será que $Y(\mathbf{x})$ apresenta um comportamento sistemático?
 - ▶ Exemplo: Seja $Y(\mathbf{x})$ a profundidade do fundo do mar, então $E[Y(\mathbf{x})]$ aumenta regularmente a partir da praia.
- O modelo puramente **intrínseco situa-se entre os casos estacionários e não estacionários**. A escolha do grau de não estacionariedade depende do caso em estudo.

8 / 26

Duas outras propriedades importantes

- $Y(x)$ é **isotrópico**, se permanece invariante quando sujeito a rotações de coordenadas (oposto a anisotrópico). Por exemplo, se o processo aleatório intrínseco $Y(x)$ é isotrópico, então $\text{Var}[Y(x_i) - Y(x_j)] = 2\gamma(\|x_i - x_j\|)$ ³



só importa a norma e não a direção

- $Y(x)$ é **ergódico**, se a média de todas as realizações possíveis for igual à média de uma única realização (permite estimativas de parâmetros com apenas 1 realização)

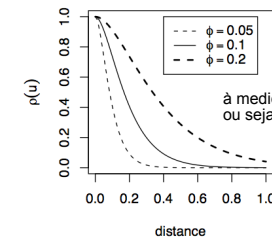
³Como $\|\cdot\|$ denota a norma Euclidiana, então γ apenas depende da distância entre as duas localizações e não da direção do vetor formado pelas duas localizações

Vantagens do variograma (Comparado com a FAC)

- O variograma adapta-se mais facilmente a observações não estacionárias (ex. $\mu(x)$ não constante)
- Para estimar o variograma, nenhuma estimativa de μ é necessária
- A estimação do variograma é mais simples que a estimação da função de covariância

As funções de covariância e variograma

- Se $\text{Var}[Y(x)] = \sigma^2$, então pode-se escrever a função de covariância⁴ como $c(u) = \sigma^2 \rho(u; \phi)$, onde $\rho(\cdot)$ é a **função de correlação**⁵



à medida que ϕ aumenta, temos correlação para maiores distancias, ou seja, temos correlação durante mais tempo

- Para $Y(x)$ sem “erro de medição”, o variograma pode ser escrito como

$$\gamma(u) = c(0) - c(u) = \sigma^2 - \sigma^2 \rho(u; \phi) = \sigma^2(1 - \rho(u; \phi)) \quad (1)$$

- O variograma $\gamma(u)$ mede a **desassociação** entre variáveis⁶

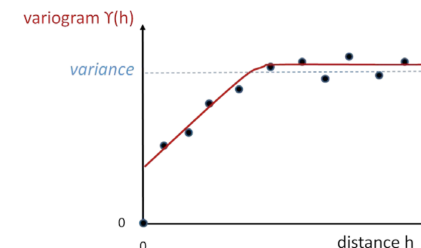
⁴ Assume-se que $Y(x)$ é um processo estacionário e isotrópico de 2ª ordem.

⁵ Parâmetro ϕ relacionado com a distância além da qual a correlação entre as variáveis é 0 (“raio de influência” ou *range*)

⁶ Oposto à função de covariância $c(u)$, que mede a **associação** entre variáveis

Análise estrutural via o variograma

- A análise estrutural tem como objetivo capturar, descrever e **modelar a maneira como uma variável geo-referenciada é estruturada espacialmente**.
- O **variograma** mede a **variabilidade média entre dois pontos quaisquer** como função do vetor de distância entre esses pontos.
 - 1 Primeiro, calculamos o variograma experimentalmente.
 - 2 Em seguida, ajustamos o variograma e modelamos a variável de interesse.



as coisas mais proximas estao mais associadas, logo as mais perto tb estao menos desassociadas

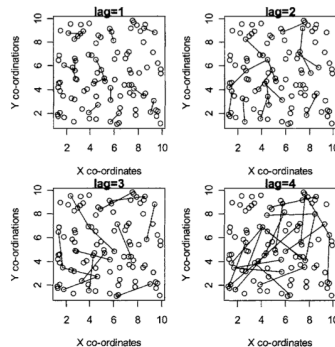
Variograma experimental $\hat{\gamma}(u)$

$$\gamma(\|x_i - x_j\|) = \frac{1}{2} \text{Var} [Y(x_i) - Y(x_j)] = \frac{1}{2} E [(Y(x_i) - Y(x_j))^2]$$

$\hat{\gamma}(u)$ pode ser obtido aproximando $E[\cdot]$ por uma média amostral

$$\hat{\gamma}(u) = \frac{1}{2|N(u)|} \sum_{N(u)} (Y(x_i) - Y(x_j))^2 \quad (2)$$

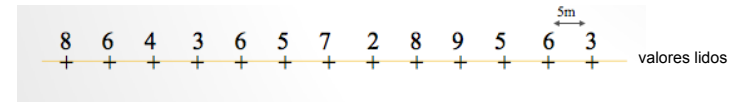
onde $N(u) = \{(x_i, x_j) : \|x_i - x_j\| \approx u, u \in \mathbb{R}\}$ e $|N(u)| = \text{"n. de pares em } N(u)\text{"}$.



13 / 26

Exercício: variograma para uma amostra regular 1-D

Variável Z definida numa grelha regular em 1-D

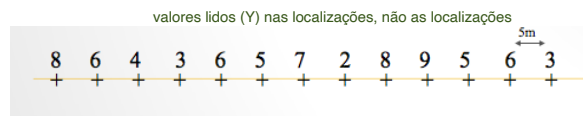


- Calcule o variograma experimental para os defasamentos: 5m, 10m e 15m.
- Avalie o variograma experimental da nova variável:

$$Y(x) = Z(x) + 3.2$$

14 / 26

Solução: variograma para uma amostra regular 1-D



- Considere a distância de 5m:

$$\hat{\gamma}(5m) = \frac{1}{2 \times 12} [(8-6)^2 + (6-4)^2 + (4-3)^2 + \dots + (6-3)^2] = 4.625$$

- Considere a distância de 10m:

$$\hat{\gamma}(10m) = \frac{1}{2 \times 11} [(8-4)^2 + (6-3)^2 + (4-6)^2 + \dots + (5-3)^2] = 5.227$$

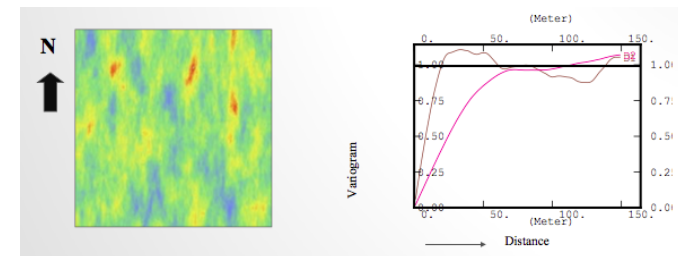
- Considere a distância de 15m:

$$\hat{\gamma}(15m) = \frac{1}{2 \times 10} [(8-3)^2 + (6-6)^2 + (4-5)^2 + \dots + (9-3)^2] = 6.000$$

Quanto **maior a distância**, menor o número de pares na linha, e **menos $\hat{\gamma}(\cdot)$ é representativo da variabilidade de Z (e de Y)**.

15 / 26

Exemplo – falha de isotropia



- Porque não devemos assumir isotropia, mas sim anisotropia?
 - Cálculo de 2 variogramas, um para a direção E-O e outro para N-S, origina resultados distintos (painel direita)
- Cada variograma direcional usa apenas pares ao longo da direção considerada, com uma **tolerância na direção**.

16 / 26

Variograma direcional versus ou omnidirecional

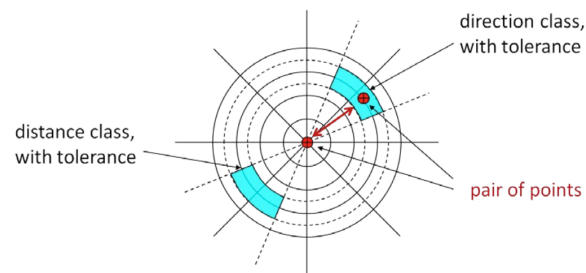


Figure: Cálculo do variograma com **tolerâncias na distância e direção**: cada par de pontos de dados é alocado a uma “categoria” de distância e direção, levando em consideração as tolerâncias.

17 / 26

Por que $\hat{\gamma}(u)$ não deve ser usado para inferência e predição?

O variograma experimental pode não ser definido-negativo condicional, o que pode levar a valores negativos absurdos para o erro de predição quadrático médio (considerado um estimador inválido).

Como obter um estimador de variograma válido?

Uma abordagem comum é aproximar $\hat{\gamma}(u)$ por algum modelo teórico $\gamma(u; \theta)$, conhecido por ser válido, capturando a dependência espacial subjacente aos dados disponíveis.

- Primeiro, escolhe-se um modelo teórico (por exemplo, exponencial, esférico, ...), normalmente usando ferramentas gráficas, que depende dos parâmetros θ
- Depois, recorre-se a:
 - ▷ um critério de ajuste clássico (por ex. mínimos quadrados) para completar a especificação de γ final, estimando os parâmetros θ , tal como se segue

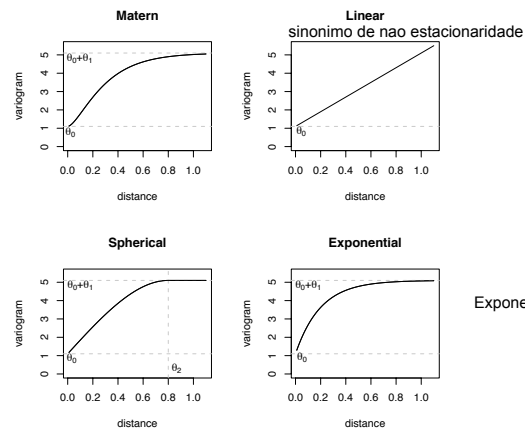
$$\hat{\theta} = \min_{\theta} \left\{ \sum_i w_i (\hat{\gamma}(u_i) - \gamma(u_i; \theta))^2 \right\}$$

- ▷ ou uma abordagem baseada em métodos de máxima verossimilhança

18 / 26

Modelos de variograma isotrópico $\gamma(u; \theta)$

Assumindo-se que $Y(x)$ tem associado um erro de medição



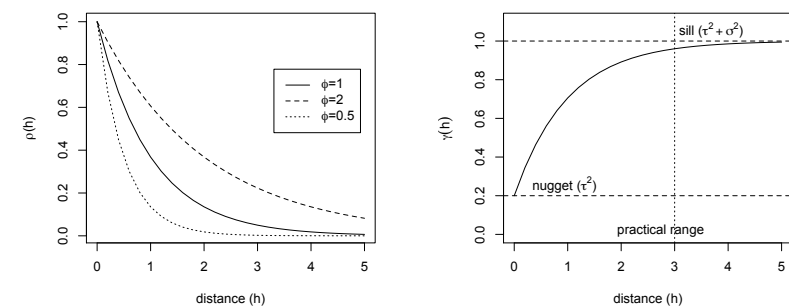
Exponencial é um caso particular da Matern, qd $k=0.5$

- ▷ $\theta_0 = \tau^2$ é conhecido como **variância do erro (nugget)**
- ▷ $\theta_0 + \theta_1 = \tau^2 + \sigma^2$ conhecido como **variância total de $Y(x)$ (sill)** "teto" onde estabiliza
- ▷ $\theta_2 = \phi$ conhecido como **raio de influência (range)**

19 / 26

A função de correlação e o variograma

Processo isotrópico e estacionário de segunda ordem



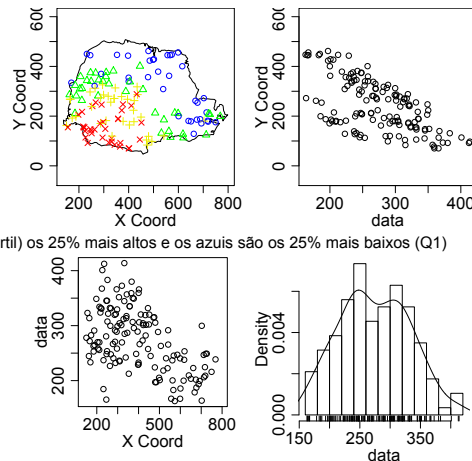
Painel da esquerda: Função de correlação exponencial para diferentes valores de range ϕ

Painel da direita: Representação esquemática de um variograma típico com seus parâmetros estruturais, sendo $\gamma(u) = \tau^2 + \sigma^2(1 - \rho(u; \phi)) = \frac{1}{2} \text{Var}[Y(x) - Y(x')]$ e $u = \|x - x'\|$

20 / 26

Exemplo 1: Dados de precipitação, Estado do Paraná, Brasil

Em muitas aplicações práticas, é útil considerar um processo Gaussiano espacial com uma função média ou tendência **possivelmente dependendo da localização, i.e. $\mu(\mathbf{x})$** , mas com uma estrutura de covariância estacionária

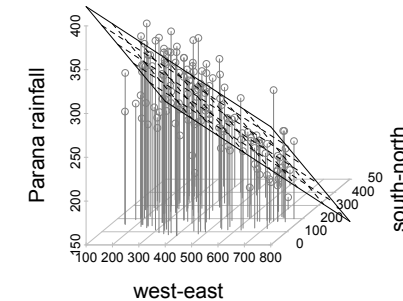


Os vermelhos são acima do Q3 (quartil) os 25% mais altos e os azuis são os 25% mais baixos (Q1)

21 / 26

Exemplo 1: Dados de precipitação - função média $\mu(\mathbf{x})$

Uma solução possível é especificar $\mu(\mathbf{x})$ como um modelo de regressão, dependendo das próprias coordenadas Fazemos isto para tentar estacionarizar a media

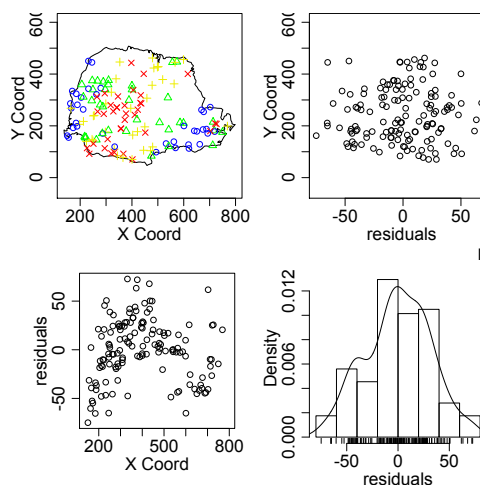


Na presença de variáveis explicativas geo-referenciadas relevantes, é razoável incluí-las em $\mu(\mathbf{x})$. Por exemplo, suponhamos que o objetivo é analisar as temperaturas médias diárias em toda a Argentina, então poderá ser importante considerar a altitude como uma covariável

22 / 26

Exemplo 1: Dados de precipitação - removendo $\mu(\mathbf{x})$

$Y(\mathbf{x}) - \mu(\mathbf{x})$ torna-se um processo estacionário Gaussiano de média zero



Neste caso já não há padrão o que é bom

23 / 26

Diferentes escalas de variabilidade

Pode-se escrever o processo aleatório espacial como

$$Y(\mathbf{x}) = \mu(\mathbf{x}) + S(\mathbf{x}) + \epsilon(\mathbf{x}) = \alpha + \sum_{j=1}^p \beta_j X_j(\mathbf{x}) + S(\mathbf{x}) + \epsilon(\mathbf{x})$$

- $\mu(\mathbf{x})$ representa uma tendência **determinística**, e identifica uma **variabilidade de grande-escala**. $\mu(\mathbf{x})$ pode envolver variáveis explicativas X_j , eventualmente dependentes da localização \mathbf{x}
- $S(\mathbf{x})$ é um processo **aleatório** com média zero e $\text{Var}[S(\mathbf{x})] = \sigma^2$, que contém a **estrutura de dependência espacial**, com $\text{Corr}[S(\mathbf{x}), S(\mathbf{x}')] = \rho(\|\mathbf{x} - \mathbf{x}'\|; \phi)$, identificando uma **variabilidade de pequena-escala**
- $\epsilon(\mathbf{x}) \sim N(0, \tau^2)$ é um **erro de medição**, i.i.d. e independente de $S(\mathbf{x})$ ⁷

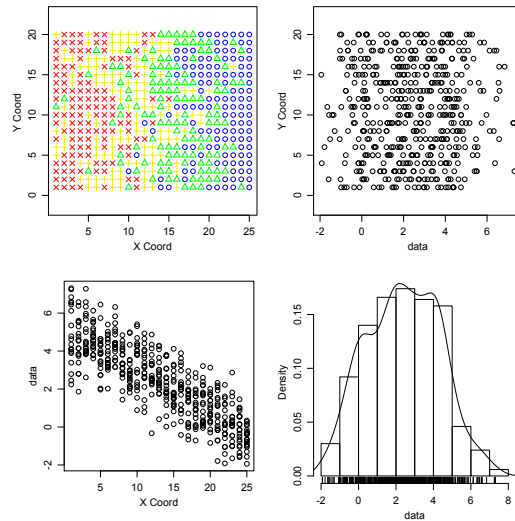
⁷A variância do erro de medição deve ser adicionada ao variograma de $Y(\mathbf{x})$ na equação (1), ficando $\gamma(u) = \tau^2 + \sigma^2(1 - \rho(u; \phi))$.

24 / 26

Exemplo 2: Dados de grãos de trigo (wheat data)



Forte tendência oeste-este (90°)

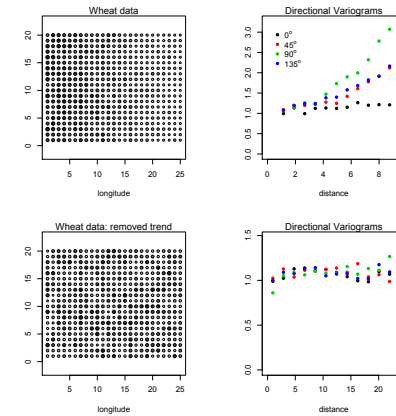


25 / 26

Exemplo 2: Trigo – detectando e removendo a tendência

$$\text{Grão}(x) = \alpha + \beta_1 \text{Longitude}(x) + \beta_2 \text{Latitude}(x) + S(x) + \epsilon(x)$$

	estimativa	s.e.	p-valor
α	5.283	0.129	<0.001
β_1	-0.226	0.006	<0.001
β_2	0.004	0.008	0.599



26 / 26