

# Máxima Verossimilhança

Prof. Wagner H. Bonat

Universidade Federal do Paraná  
Departamento de Estatística  
Laboratório de Estatística e Geoinformação



# Sumário

1 Notação e definições

2 Verossimilhança

# Notação e Definições

- O vetor ( $n \times 1$ ) de variáveis aleatórias (va) é denotado por  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ .
- O vetor ( $n \times 1$ ) de realizações de uma va é denotado por  $\mathbf{y} = (y_1, \dots, y_n)^\top$ .
- Denote  $f(\mathbf{Y}; \theta)$  a função de probabilidade (fp) caso discreto ou função de densidade probabilidade (fdp) do vetor aleatório  $\mathbf{Y}$ .
- **Parâmetro ou vetor de parâmetros** Vetor de características numéricas de uma população.
- Notação:  $\theta$  vetor ( $p \times 1$ ) de parâmetros.
- **Espaço paramétrico** é o conjunto de todas as possíveis combinações entre todos os valores para todos os diferentes parâmetros envolvidos em uma fp ou fdp. Notação  $\Theta$ .
- **Suporte** é conjunto de valores realizáveis de uma va. Notação  $\Omega$ .
- Exemplos: Binomial, Poisson e normal.

# Definições

- Uma **estatística** é uma variável aleatória  $T = t(\mathbf{Y})$ , onde a função  $t(\cdot)$  não depende de  $\theta$ .
- Uma **estatística**  $T$  é um **estimador** para  $\theta$  se o valor realizado  $t = t(\mathbf{y})$  é usado como uma **estimativa** para o valor de  $\theta$ .
- A distribuição de probabilidade de  $T$  é chamada de **distribuição amostral** do estimador  $t(\mathbf{Y})$ .
- O **viés** de um estimador  $T$  é a quantidade

$$B(T) = E(T - \theta).$$

O estimador  $T$  é dito não viciado para  $\theta$  se  $B(T) = 0$ , tal que  $E(T) = \theta$ .

- O estimador  $T$  é assintoticamente não viciado para  $\theta$  se  $E(T) \rightarrow \theta$  quando  $n \rightarrow \infty$ .

# Definições

- A **eficiência relativa** entre dois estimadores  $T_1$  e  $T_2$  é a razão  $er = \frac{V(T_1)}{V(T_2)}$  em que  $V(\cdot)$  denota a variância.
- O **erro quadrático médio** de um estimador  $T$  é a quantidade

$$EQM(T) = E((T - \theta)^2) = V(T) + B(T)^2.$$

- Um estimador  $T$  é **médio quadrático consistente** para  $\theta$  se o  $EQM(T) \rightarrow 0$  quando  $n \rightarrow \infty$ .
- O estimador  $T$  é **consistente em probabilidade** se  $\forall \epsilon > 0$ ,  $P(|T - \theta| > \epsilon) \rightarrow 0$ , quando  $n \rightarrow \infty$ .

# Sumário

1 Notação e definições

2 Verossimilhança

# Função de verossimilhança

- Sejam dados  $\mathbf{y}$  uma realização de um vetor aleatório  $\mathbf{Y}$  com fp ou fdp  $f(\mathbf{Y}, \theta)$ . A **função de verossimilhança**  $L(\theta) : \Theta \rightarrow [0, \infty]$  para  $\theta$  é a função aleatória

$$L(\theta) \equiv f(\mathbf{Y}, \theta)$$

onde  $f(y_1, \dots, y_n | \theta)$  é a função de distribuição conjunta de  $\mathbf{Y}$ .

- 1 Caso discreto não há ambiguidade então

$$L(\theta) \equiv P_\theta[Y = y].$$

- 2 Caso contínuo em geral as observações são medidas com algum grau de precisão em um intervalo ( $y_{iI} \leq y_i \leq y_{iS}$ ). Neste caso a verossimilhança é dada por

$$L(\theta) = P_\theta[y_{1I} \leq y_1 \leq y_{1S}, y_{2I} \leq y_2 \leq y_{2S}, \dots, y_{nI} \leq y_n \leq y_{nS}].$$

# Função de verossimilhança

- Suponha que as observações são independentes e medidas com o mesmo grau de precisão.
- Assim, cada dado é medido em um intervalo  $(y_i - \delta/2 \leq Y_i \leq y_i + \delta/2)$ .
- Com estas suposições a verossimilhança pode ser escrita como

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n P_{\theta}[y_i - \delta/2 \leq Y_i \leq y_i + \delta/2] \\ &= \prod_{i=1}^n \int_{y_i - \delta/2}^{y_i + \delta/2} f(y_i, \theta) d(y_i). \end{aligned}$$

# Função de verossimilhança

- Se o grau de precisão é alto ( $\delta$  é pequeno) em relação a variabilidade dos dados, a expressão se reduz a

$$L(\boldsymbol{\theta}) \approx \left( \prod_{i=1}^n f(y_i, \boldsymbol{\theta}) \right) \delta^n.$$

- Finalmente, se  $\delta$  não depende de  $\boldsymbol{\theta}$ , temos

$$L(\boldsymbol{\theta}) \approx \prod_{i=1}^n f(y_i, \boldsymbol{\theta}),$$

- Para enfatizar que a verossimilhança é avaliada nas observações usamos a notação  $L(\boldsymbol{\theta}|\mathbf{y})$ .

# Função de Log-Verossimilhança

- A função de log-verossimilhança é a função estocástica  $l(\boldsymbol{\theta}) : \Theta \rightarrow \mathfrak{R}$  definida por

$$l(\boldsymbol{\theta}|\mathbf{y}) = \log(L(\boldsymbol{\theta}|\mathbf{y})).$$

- No caso iid, tem-se

$$l(\boldsymbol{\theta}|\mathbf{y}) = \sum_{i=1}^n \log(L(\boldsymbol{\theta}|y_i)).$$

- $l(\boldsymbol{\theta}|\mathbf{y}) = -\infty$  quando  $L(\boldsymbol{\theta}) = 0$ , mas isso ocorre quando  $f(y_1, \dots, y_n|\boldsymbol{\theta}) = 0$  que tem probabilidade de ocorrência igual a zero.

# Vetor escore

- O vetor escore  $\mathbf{U}(\boldsymbol{\theta}|\mathbf{y}) : \Theta \rightarrow \mathbb{R}^p$  é um vetor aleatório  $p \times 1$  definido por

$$\mathbf{U}(\boldsymbol{\theta}|\mathbf{y}) = \frac{\partial l(\boldsymbol{\theta}|\mathbf{y})}{\partial \boldsymbol{\theta}} = \begin{pmatrix} \frac{\partial l(\boldsymbol{\theta}|\mathbf{y})}{\partial \theta_1} \\ \vdots \\ \frac{\partial l(\boldsymbol{\theta}|\mathbf{y})}{\partial \theta_p} \end{pmatrix}.$$

- Notação popular em termos de gradiente

$$\mathbf{U}(\boldsymbol{\theta}|\mathbf{y}) = \nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta}|\mathbf{y}).$$

- Notação  $j$ -ésimo componente de  $\mathbf{U}(\boldsymbol{\theta}|\mathbf{y})$  por  $U_j(\boldsymbol{\theta}|\mathbf{y})$ .

# Matriz de informação esperada (Fisher Information Matrix)

- A matriz  $p \times p$  definida por

$$\begin{aligned}\mathbf{I}_E(\boldsymbol{\theta}) &= \mathbf{Var}(\mathbf{U}(\boldsymbol{\theta}|\mathbf{Y})) \\ &= \mathbf{E}(\mathbf{U}(\boldsymbol{\theta}|\mathbf{Y})\mathbf{U}^\top(\boldsymbol{\theta}|\mathbf{Y})).\end{aligned}$$

é chamada de matriz de informação esperada.

- As entradas  $j$  e  $k$  são expressadas por

$$\begin{aligned}\mathbf{I}_{Ejk}(\boldsymbol{\theta}) &= \text{Cov}(\mathbf{U}_j(\boldsymbol{\theta}|\mathbf{Y}), \mathbf{U}_k(\boldsymbol{\theta}|\mathbf{Y})) \\ &= \mathbf{E}(\mathbf{U}_j(\boldsymbol{\theta}|\mathbf{Y})\mathbf{U}_k(\boldsymbol{\theta}|\mathbf{Y})).\end{aligned}$$

# Matriz de informação observada

- A matriz  $p \times p$  definida por

$$I_O(\boldsymbol{\theta}) = -\frac{\partial^2 l(\boldsymbol{\theta}|\mathbf{y})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}$$

é chamada de informação observada.

- As entradas  $j$  e  $k$  da matriz de informação observada é dada por

$$I_{Ojk}(\boldsymbol{\theta}) = -\frac{\partial^2 l(\boldsymbol{\theta}|\mathbf{y})}{\partial \theta_j \partial \theta_k}.$$

- Primeira igualdade de Bartlett  $E(\mathbf{U}(\boldsymbol{\theta}|\mathbf{Y})) = \mathbf{0}$ .
- Segunda igualdade de Bartlett

$$I_E(\boldsymbol{\theta}) = E(I_O(\boldsymbol{\theta})).$$

# Estimador de máxima verossimilhança (EMV)

- Estimativa de máxima verossimilhança: Seja  $L(\boldsymbol{\theta}, \mathbf{y})$  a função de verossimilhança. O valor  $\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}(\mathbf{y})$  é a estimativa de máxima verossimilhança para  $\boldsymbol{\theta}$  se  $L(\hat{\boldsymbol{\theta}}) \geq L(\boldsymbol{\theta}), \forall \boldsymbol{\theta}$ .
- Estimador de máxima verossimilhança: Se  $\hat{\boldsymbol{\theta}}(\mathbf{y})$  é a estimativa de máxima verossimilhança, então  $\hat{\boldsymbol{\theta}}(\mathbf{Y})$  é o estimador de máxima verossimilhança (EMV).
- Em geral,  $\hat{\boldsymbol{\theta}}$  satisfaz a equação de verossimilhança

$$U(\boldsymbol{\theta}|\mathbf{y}) = \mathbf{0}.$$

- Um sistema com  $p$  equações e  $p$  incógnitas

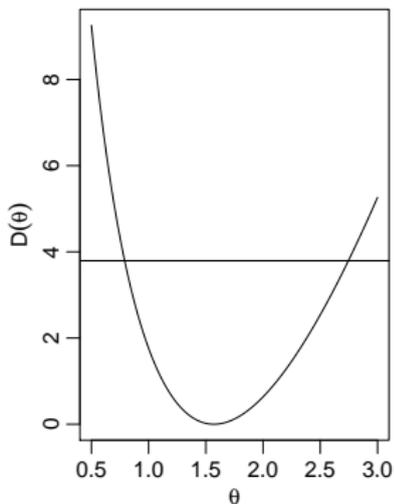
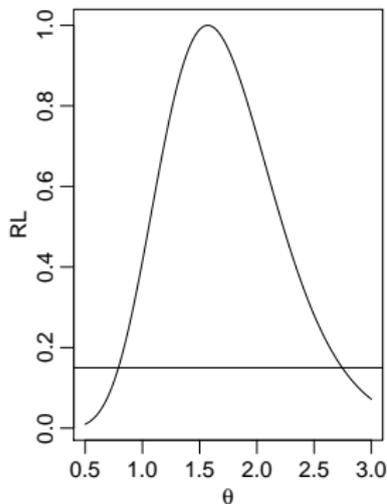
$$\begin{pmatrix} U_1(\boldsymbol{\theta}|\mathbf{y}) = 0 \\ \vdots \\ U_p(\boldsymbol{\theta}|\mathbf{y}) = 0 \end{pmatrix}$$

# Intervalos de confiança

- Um intervalo de verossimilhança para  $\theta$  é um intervalo da forma  $\theta : L(\theta) \geq rL(\hat{\theta})$  ou equivalentemente,  $\theta : D(\theta) \leq c^*$ , com  $D(\theta) = -2[l(\theta) - l(\hat{\theta})]$  e  $c^* = -2 \log(r)$ .
- Para o caso multiparâmetros a ideia se mantém, mas trocamos o intervalo por uma região de confiança.
- Nesta definição  $r$  precisa ser um valor entre 0 e 1, logo  $c^* > 0$  para intervalos não-vazios.

## Estratégias para construir o intervalo de confiança

- Verossimilhança relativa  $\frac{L(\theta)}{L(\hat{\theta})} \geq r$ .
- Deviance  $D(\theta) = -2[l(\theta) - l(\hat{\theta})] \leq -2 \log(r)$ .
- Após definir o valor  $c^* = -2 \log(r)$ , é necessário encontrar as raízes da função de verossimilhança relativa ou da *deviance*.



## Estratégias para construir o intervalo de confiança

- De forma geral métodos numéricos serão necessários.
- Alternativa: Expansão em séries de Taylor (segunda ordem) para  $l(\theta)$  em torno de  $\hat{\theta}$

$$D(\theta) = -2[l(\theta) - l(\hat{\theta})] \\ \approx 2 \left\{ l(\hat{\theta}) - [l(\hat{\theta}) + (\theta - \hat{\theta})l'(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^2 l''(\hat{\theta})] \right\}.$$

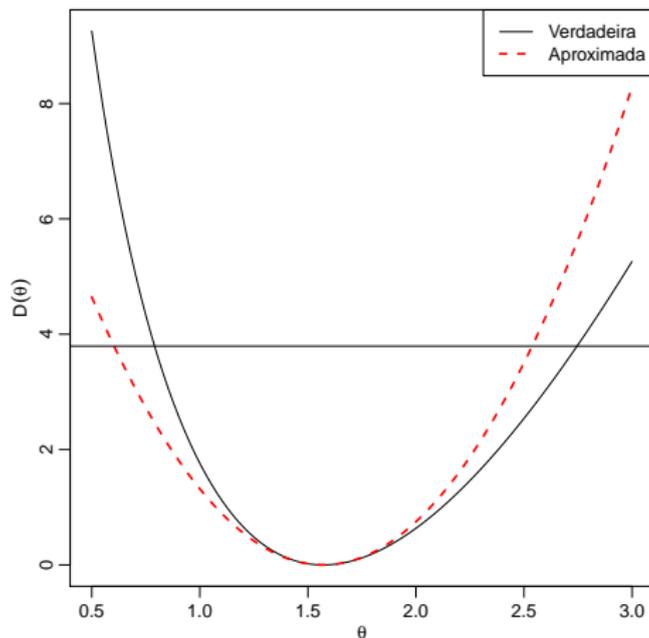
- Eliminando termos, temos

$$D(\theta) \approx -(\theta - \hat{\theta})^2 l''(\hat{\theta}) \leq c^*.$$

- Resolvendo em  $\theta$  chegamos a intervalos da forma.

$$\hat{\theta} \pm \sqrt{\frac{c^*}{-l''(\hat{\theta})}}.$$

## Estratégias para construir o intervalo de confiança



- Ainda precisamos definir  $r$  ou  $c^*$ .

# Condições de regularidade

- O parâmetro  $\theta$  é **identificável**. Isso significa que se  $f(\theta_1|\mathbf{y}) = f(\theta_2|\mathbf{y})$  para quase todos  $\mathbf{y} \in \mathfrak{R}$ , então  $\theta_1 = \theta_2$ .
- O suporte de  $f(\theta|\mathbf{y})$  é o mesmo para todo  $\theta \in \mathfrak{R}$ .
- O verdadeiro valor do parâmetro  $\theta_0$  pertence ao interior de  $\Theta$ .
- $f(\theta|\mathbf{y})$  é duas vezes continuamente diferenciável com relação aos componentes de  $\theta$  para quase todo  $\mathbf{y} \in \mathfrak{R}$ .
- $\frac{\partial}{\partial \theta}$  e  $\int$  (caso contínuo), ou  $\frac{\partial}{\partial \theta}$  e  $\sum$  (caso discreto) podem ser intercambiados.

# Propriedades do Estimador de máxima verossimilhança

- Sendo  $\theta_0$  o verdadeiro valor do vetor de parâmetros  $\theta$ . Então
- Consistência:  $\hat{\theta} \xrightarrow{P} \theta_0$ , ou seja,

$$P\left(\|\hat{\theta} - \theta_0\| > \epsilon\right) \rightarrow 0, \quad \text{quando } n \rightarrow \infty.$$

- Normalidade assintótica:

$$\hat{\theta} \xrightarrow{D} N_p(\theta, I_E^{-1}(\theta)), \quad \text{quando } n \rightarrow \infty.$$

- Qualquer termo assintoticamente equivalente a  $I_E(\theta)$  pode ser usado.

$$\hat{\theta} \sim NM_p(\theta, I_E^{-1}(\hat{\theta}))$$

$$\hat{\theta} \sim NM_p(\theta, I_O^{-1}(\theta))$$

$$\hat{\theta} \sim NM_p(\theta, I_O^{-1}(\hat{\theta})).$$

# Distribuição assintótica da deviance

- Para um problema regular de estimação, no limite com  $n \rightarrow \infty$ , se  $\theta$  é o verdadeiro valor do parâmetro, então

$$D(\theta) = -2[l(\theta) - l(\hat{\theta})] \sim \chi_p^2$$

ou seja, a função deviance segue uma distribuição Qui-Quadrado com  $d$  graus de liberdade, onde  $p$  é a dimensão do vetor  $\theta$ .

## Resumo dos resultados

- O estimador de máxima verossimilhança  $\hat{\theta}$  de  $\theta$  é assintoticamente não-viciado, isto é,  $E(\hat{\theta}) \rightarrow \theta$ .
- Assintoticamente  $V(\hat{\theta}) \rightarrow I_E^{-1}(\theta)$ , o qual por uma versão multivariada do limite de Cramér-Rao é o melhor possível, mostrando que o EMV é eficiente para o vetor  $\theta$ .
- Denote  $J = I_E^{-1}(\theta)$ , então  $V(\hat{\theta}) = J$ , sendo que,  $J$  é uma matriz simétrica e definida positiva, com elementos  $J_{ij} = \text{Cov}(\hat{\theta}_i, \hat{\theta}_j)$  então  $J_{ii}$  é a variância de  $\hat{\theta}_i$ .
- Denota-se  $J_{ii}^{\frac{1}{2}}$  de desvio padrão de  $\hat{\theta}_i$ .
- Podemos construir intervalos de  $100(1 - \alpha)\%$  de confiança para  $\theta_i$  na forma  $\hat{\theta}_i \pm z_{\frac{\alpha}{2}} J_{ii}^{\frac{1}{2}}$ . Intervalos desta forma serão denominados, intervalos de Wald ou baseados em aproximação quadrática da verossimilhança.

## Regiões de confiança

- Para regiões de confiança baseados na *deviance* considera-se  $\theta \in \Theta : D(\theta) \leq c^*$ , para algum valor  $c^*$  a ser especificado. Pode-se escolher  $c^*$  baseado em justificativas assintóticas de que  $D(\theta) \sim \chi_p^2$  é uma escolha razoável para  $c^* = c_\alpha$  com  $P(\chi_p^2 \geq c_\alpha) = \alpha$ , por exemplo se  $\alpha = 0.05$  e  $p = 1$ , então  $c_\alpha = 3.84$ . Isto gera uma região de  $100(1 - \alpha)\%$  de confiança. Estes intervalos serão denominados de intervalos *deviance*.
- Note que intervalos *deviance* e baseados em verossimilhança são de certa forma equivalentes. A diferença é a justificativa para a forma da região de confiança.
- Discutir os tutoriais I e II.

# Método delta

- Considere obter um intervalo de confiança para  $\phi = g(\boldsymbol{\theta})$  por invariância temos que  $\hat{\phi} = g(\hat{\boldsymbol{\theta}})$  e a variância de  $\hat{\phi}$  é dada por

$$V(\hat{\phi}) = V(g(\hat{\boldsymbol{\theta}})) = \nabla g(\hat{\boldsymbol{\theta}})^\top I_E(\hat{\boldsymbol{\theta}})^{-1} \nabla g(\hat{\boldsymbol{\theta}})$$

onde

$$\nabla g(\hat{\boldsymbol{\theta}}) = \left( \frac{\partial g(\hat{\boldsymbol{\theta}})}{\partial \theta_1}, \dots, \frac{\partial g(\hat{\boldsymbol{\theta}})}{\partial \theta_d} \right)^\top.$$

- De imediato, temos

$$\hat{\phi} \sim N(\phi, \nabla g(\boldsymbol{\theta})^\top I_E(\boldsymbol{\theta})^{-1} \nabla g(\boldsymbol{\theta})).$$

- Ver tutorial III.

# Verossimilhança perfilhada

- A verossimilhança perfilhada de  $\theta$  é definida por

$$\tilde{L}(\theta) = L(\theta, \hat{\lambda}_\theta).$$

- A forma apresentada sugere um procedimento de maximização em duas etapas. A primeira consiste em obter  $\hat{\lambda}_\theta$  que maximiza  $l(\theta, \lambda) = \log L(\theta, \lambda)$  com respeito a  $\lambda$  supondo  $\theta$  fixo.
- A seguir maximiza-se  $\tilde{l}(\theta)$ . Assim, uma região ou intervalo de confiança para  $\theta$  pode ser obtida usando que

$$\tilde{D}(\theta) = -2[\tilde{l}(\theta) - \tilde{l}(\hat{\theta})] \sim \chi_p^2$$

onde  $p$  é a dimensão de  $\theta$ .

# Parâmetros ortogonais

- Considere um modelo estatístico parametrizado por  $\theta = (\theta_1, \theta_2)^\top$ . No caso da matriz de informação de Fisher ser bloco diagonal

$$I_E(\theta) = \begin{pmatrix} I_{E_1}(\theta) & 0 \\ 0 & I_{E_2}(\theta) \end{pmatrix},$$

os vetores de parâmetros  $\theta_1$  e  $\theta_2$  são ditos **ortogonais**.

- A distribuição assintótica de  $\theta_1$  é a mesma se  $\theta_2$  é considerado conhecido ou desconhecido.
- Definição similar pode ser feita usando a matriz de informação observada.
- Ver tutorial II.