

# Inferência Estatística com Abordagem Bayesiana

Rosangela Helena Loschi <sup>1</sup>

<sup>1</sup>Departamento de Estatística  
Universidade Federal de Minas Gerais

16 de dezembro de 2021

# INFERÊNCIA INTERVALAR

- ▶ Regiões de credibilidade
- ▶ Regiões com mais alta densidade *a posteriori*

# Estimação por Regiões

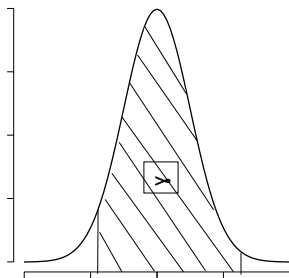
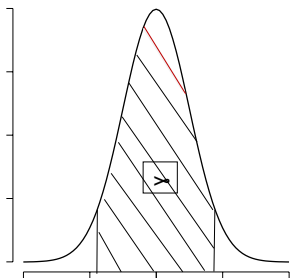
- ▶ As estimativas pontuais são resumos muito radicais da distribuição *a posteriori*  $\pi(\theta | x)$ .
- ▶ As **regiões de credibilidade** são mais informativas pois são escolhidas de forma a conter parte substancial da massa probabilística da distribuição *a posteriori*.
- ▶ Uma região do espaço paramétrico  $R_\theta(x) \subseteq \Theta$  é uma região com credibilidade  $\gamma$  se

$$P(\theta \in R_\theta(x) | x) \geq \gamma \quad (1)$$

- ▶ Obs: No caso contínuo

$$P(\theta \in R_\theta(x) | x) = \int_{R_\theta(x)} \pi(\theta | x) d\theta = \gamma$$

# Estimação por Regiões: Regiões HPD

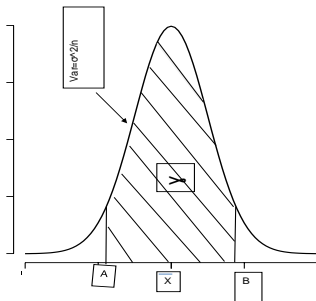


- ▶ Uma região de credibilidade tem uma interpretação probabilística direta pois não é aleatória.
- ▶ Concluimos que a probabilidade *a posteriori* de  $\theta$  pertencer à região de credibilidade  $R_\theta(x)$  é  $\gamma$ .

## Estimação por Regiões

**Exemplo:** Suponha que  $X|\mu \stackrel{iid}{\sim} \text{Normal}(\mu, \sigma^2)$ ,  $\sigma^2$  conhecido, e que a incerteza *a priori* sobre  $\mu$  é bem descrita pela distribuição de Jeffreys. Encontre regiões de  $\mathbb{R}$  com credibilidade  $\gamma$ .

- ▶ A distribuição *a posteriori* é  $\mu | x \sim \text{Normal}(\bar{x}, \sigma^2/n)$ .
  - + simétrica e unimodal



- ▶ A região com credibilidade  $\gamma$  é um intervalo  $[A, B]$  definido sob a distribuição *a posteriori* de  $\mu$  tal que  $P(A < \mu < B | \mathbf{x}) = \gamma$

# Estimação por Regiões

- ▶ Devemos ter

$$Z_{\alpha_1} = \frac{A - \bar{x}}{\sqrt{\sigma^2/n}} \quad Z_{\alpha_2} = \frac{B - \bar{x}}{\sqrt{\sigma^2/n}}$$

- ▶ Um intervalo com credibilidade  $\gamma$  para  $\mu$  tem a seguinte forma

$$IC_{\gamma}(\mu) = \left[ \bar{x} + \frac{\sigma}{\sqrt{n}} Z_{\alpha_1}; \bar{x} + \frac{\sigma}{\sqrt{n}} Z_{\alpha_2} \right]$$

onde  $Z_{\alpha_i}$  é encontrado na tabela da distribuição  $N(0, 1)$  e tal que  $P(Z > |Z_{\alpha_i}|) = \alpha_i$ ,  $i = 1, 2$  e  $\alpha_1 + \alpha_2 = 1 - \gamma$ .

- ▶ Note que os limites são números determinados a partir da distribuição *a posteriori* para  $\mu$ .

# Estimação por Regiões

- Um intervalo com credibilidade  $\gamma$  e simétrico em torno da média *a posteriori*  $\bar{x}$  é

$$IC_{\gamma}(\mu) = \left[ \bar{x} - \frac{\sigma}{\sqrt{n}} Z_{\frac{1-\gamma}{2}}; \bar{x} + \frac{\sigma}{\sqrt{n}} Z_{\frac{1-\gamma}{2}} \right]$$

onde  $Z_{\frac{1-\gamma}{2}}$  é encontrado na tabela da distribuição  $N(0, 1)$  e tal que  $P(Z > Z_{\frac{1-\gamma}{2}}) = \frac{1-\gamma}{2}$ .

## Estimação por Regiões

**Coloquemos números:** Assuma que  $\bar{x} = 10$ ,  $\sigma^2 = 4$  e  $n = 4$ .  
Construa intervalos com credibilidade  $\gamma = 0,95$  e compare-os.

- ▶ Se  $\alpha_1 = \alpha_2 = 0,025$

$$ICr(\mu \mid \mathbf{x}) = [10 - 1,96; 10 + 1,96] = [8,04; 11,96] \Rightarrow Amp = 3,92.$$

$\Rightarrow$  A probabilidade da média populacional estar entre 8,04 e 11,96 é 0,95.

- ▶ Se  $\alpha_1 = 0,03$  e  $\alpha_2 = 0,02$

$$ICr(\mu \mid \mathbf{x}) = [10 - 1,88; 10 + 2,05] = [8,12; 12,05] \Rightarrow Amp = 3,93$$

- ▶ Se  $\alpha_1 = 0,01$  e  $\alpha_2 = 0,04$

$$ICr(\mu \mid \mathbf{x}) = [10 - 2,33; 10 + 1,75] = [7,67; 11,75] \Rightarrow Amp = 4,08$$

- ▶ Se  $\alpha_1 = 0,00$  e  $\alpha_2 = 0,05$

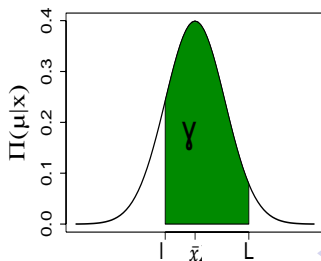
$$ICr(\mu \mid \mathbf{x}) = [-\infty; 10 + 1,64] = [-\infty; 11,64] \Rightarrow Amp = \infty$$

Todos são intervalos com credibilidade 0,95 e todos são interpretados probabilisticamente. Como escolher?



## Estimação por Regiões

- ▶ As regiões de credibilidade como definida na equação (1) são, em certo sentido, melhores resumos *a posteriori* do que as estimativas pontuais.
- ▶ No entanto, existem infinitas (no caso contínuo) regiões com credibilidade  $\gamma$ .
- ▶ Como está definida em (1), a região de credibilidade pode conter valores de  $\theta$  que estejam na regiões com menos massa probabilística que alguns outros  $\theta$  que estejam fora da região construída.



# Estimação por Regiões

- ▶ Buscamos a **melhor**! Como definir melhor?
- ▶ Buscamos a região que contém todos os valores de  $\theta$  mais prováveis *a posteriori*.

# Estimação por Regiões: regiões HPD

Uma região do espaço paramétrico  $R_\theta(x) \subseteq \Theta$ , com credibilidade  $\gamma$ , é uma **região com mais alta densidade *a posteriori*** (Highest Posterior Density) se

(1)

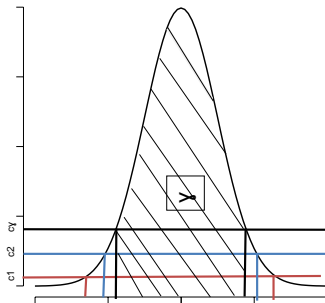
$$R_\theta(x) = \{\theta \in \Theta : \pi(\theta | x) \geq C_\gamma\}$$

onde  $C_\gamma > 0$  é a maior constante tal que

(2)

$$P(\theta \in R_\theta(x) | x) \geq \gamma.$$

# Estimação por Regiões: Regiões HPD



- ▶ Se a região HPD é um intervalo  $[A, B] \rightarrow \pi(A | x) = \pi(B | x)$ .
- ▶ A região HPD é única.
- ▶ O intervalo HPD é o de menor amplitude (A região HPD é a de menor área ou menor volume).

## Estimação por Regiões; Regiões HPD

**Exemplo: (cont.)** Suponha que  $X|\mu \stackrel{iid}{\sim} \text{Normal}(\mu, \sigma^2)$ ,  $\sigma^2$  conhecido, e que a incerteza *a priori* sobre  $\mu$  é bem descrita pela distribuição de Jeffreys. Encontre o intervalo HPD com credibilidade  $\gamma$ .

**Solução:** Como a distribuição *a posteriori* é  $\mu | \mathbf{x} \sim N(\bar{x}, \sigma^2/n)$  é unimodal e simétrica em torno de  $\bar{x}$ , o intervalo simétrico em torno da média *a posteriori*  $\bar{x}$  o Intervalo HPD

$$IC_{\gamma}(\mu) = \left[ \bar{x} - \frac{\sigma}{\sqrt{n}} Z_{\frac{1-\gamma}{2}}; \bar{x} + \frac{\sigma}{\sqrt{n}} Z_{\frac{1-\gamma}{2}} \right]$$

onde  $Z_{\frac{1-\gamma}{2}}$  é encontrado na tabela da distribuição  $N(0, 1)$  e tal que  $P(Z > Z_{\frac{1-\gamma}{2}}) = \frac{1-\gamma}{2}$ .

- ▶ Não foi coincidência que na análise numérica que fizemos, ele foi o mais curto.

# Estimação por Regiões: regiões HPD

## OBSERVAÇÕES Importantes Referentes ao Exemplo:

- ▶ Pós-experimentalmente, o intervalo clássico com confiança  $\gamma$  tem os mesmos limites que o intervalo simétrico (HPD) em torno da média *a posteriori*  $\bar{x}$ .
- ▶ Apesar disto, as inferências sobre  $\mu$  não são as mesmas.

Num contexto geral,

- ▶ A interpretação clássica é em termos de **confiança**, não é em termos de **probabilidade**;
- ▶ Em sua construção, o intervalo de confiança leva em conta todas as amostras (observadas e não) e não apenas a que foi observada.
- ▶ O intervalo HPD segue o **Princípio da Verossimilhança**.

## Estimação por Regiões: regiões HPD

**Outro Exemplo:** (Paulino, Turkeman, Murteira, 2003)

Se  $X_1, \dots, X_n | \theta \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$  e, *a priori*,  $\theta \sim \text{Beta}(2, 1)$  a distribuição *a posteriori* é  $\theta | X_1, \dots, X_n \sim \text{Beta}(9, 3)$ .

- ▶ A distribuição *a posteriori* é unimodal e assimétrica
- ▶ A região HPD será um intervalo.

Os intervalos HPD e central são

IC	90%	Amp.	95%	Amp.
HPD	[0.566;0.944]	0.378	[0.516;0.959]	0.443
Central	[0.530;0.921]	0.391	[0.482;0.94]	0.458

## Estimação por Regiões: regiões HPD

Observação importante sobre as regiões HPD

- ▶ As regiões de credibilidade (HPD ou não), se transformadas, preservam a mesma credibilidade.
- ▶ Mas as regiões HPD **não são invariantes** sob transformações um-a-um.
- ▶ Assim, se  $R_\theta(x)$  é uma região HPD para  $\theta$  e  $\phi = g(\theta)$ , a região transformada  $g(R_\theta(x))$ 
  - + não é, necessariamente, uma região HPD para  $\phi$ ;
  - + será uma região com a mesma credibilidade  $\gamma$ .

**Exemplo:** Se  $X|\theta \stackrel{iid}{\sim} \exp(\theta)$  e, *a priori*,  $\theta \sim \text{Gamma}(a, b)$  então

- + A distribuição *a posteriori* é

$$\theta | x \sim \text{Gamma}\left(n + a, b + \sum x_i\right) = \text{Gamma}(A, B).$$

$$\pi(\theta | x) \propto \theta^{A-1} \exp\{-B\theta\}$$

- + Como é uma distribuição unimodal, a região HPD é um intervalo. Denote-o por  $[\bar{\theta}; \hat{\theta}] \Rightarrow \pi(\bar{\theta} | x) = \pi(\hat{\theta} | x)$ .



## Estimação por Regiões: regiões HPD

Considere a transformação  $\alpha = 2B\theta$ . Consequentemente, segue que

- \*  $\alpha \mid \mathbf{x} \sim \text{Gamma}(A, 1/2) \Rightarrow \pi(\alpha \mid \mathbf{x}) \propto \alpha^{(A-1)} \exp\{-\alpha/2\}$
- \*  $\gamma = P(\bar{\theta} < \theta < \hat{\theta} \mid \mathbf{x}) = P(2B\bar{\theta} < \alpha < 2B\hat{\theta} \mid \mathbf{x})$
- \* Mais ainda, segue que avaliando  $\pi(\alpha \mid \mathbf{x})$  em  $2B\bar{\theta}$

$$\begin{aligned}\pi(2B\bar{\theta} \mid \mathbf{x}) &\propto (2B\bar{\theta})^{A-1} \exp\{-2B\bar{\theta}/2\} \\ &\propto (2B)^{A-1} \bar{\theta}^{A-1} \exp\{-B\bar{\theta}\} \\ &\propto (2B)^{A-1} \hat{\theta}^{A-1} \exp\{-B\hat{\theta}\} = \pi(2B\hat{\theta} \mid \mathbf{x})\end{aligned}$$

- \* Assim, o intervalo  $[2B\bar{\theta}; 2B\hat{\theta}]$  é um intervalo HPD para  $\alpha$ .

## Estimação por Regiões: regiões HPD

Considere uma outra transformação  $\alpha = \frac{1}{\theta}$ . Segue que

- \*  $\alpha \mid \mathbf{x} \sim \text{GammaInv}(A, B) \Rightarrow \pi(\alpha \mid \mathbf{x}) \propto \left(\frac{1}{\alpha}\right)^{A+1} \exp\{-B/\alpha\}$
- \*  $\gamma = P(\bar{\theta} < \theta < \hat{\theta} \mid \mathbf{x}) = P\left(\frac{1}{\hat{\theta}} < \alpha < \frac{1}{\bar{\theta}} \mid \mathbf{x}\right)$
- \* Neste caso, avaliando  $\pi(\alpha \mid \mathbf{x})$  em  $1/\bar{\theta}$  temos

$$\begin{aligned}\pi\left(\frac{1}{\bar{\theta}} \mid \mathbf{x}\right) &\propto (\bar{\theta})^2 (\bar{\theta})^{A-1} \exp\{-B\bar{\theta}\} \\ &\propto (\bar{\theta})^2 \hat{\theta}^{A-1} \exp\{-B\hat{\theta}\} \neq \pi\left(\frac{1}{\hat{\theta}} \mid \mathbf{x}\right)\end{aligned}$$

- \* Assim  $[1/\hat{\theta}; 1/\bar{\theta}]$  não é um intervalo HPD para  $\alpha$ .

# Estimação por Regiões e Teoria de decisão

- ▶ A construção de uma região HPD é um problema de decisão.
- ▶ Buscamos uma região  $R \subseteq \Theta$  com credibilidade  $\gamma$  tal que

$$\mathbb{P}(\theta \in R \mid x) = \int_R \pi(\theta \mid \mathbf{x}) d\theta = \gamma.$$

- ▶ Fixado  $\gamma$  NÃO existe apenas uma região  $R$  satisfazendo esta condição
- ▶ Nossa tarefa é buscar pela região que melhor resume a informação *a posteriori*. ← um problema de decisão.

# Estimação por Regiões e Teoria de decisão

**Proposição:** Seja  $\pi(\theta \mid \mathbf{x})$  a distribuição *a posteriori* para  $\theta$ . Seja um número real  $\gamma$ ,  $0 < \gamma < 1$  e seja o conjunto não-vazio  $\mathcal{A}$  de regiões com credibilidade  $\gamma$ , isto é,

$$\mathcal{A} = \{R \subseteq \Theta : \mathbb{P}(\theta \in R \mid \mathbf{x}) = \gamma\}$$

e considere a função de perda

$$L(R, \theta) = K\|R\|I(\theta \in R)$$

em que  $R \in \mathcal{A}$ ,  $\theta \in \Theta$ ,  $K > 0$  e  $\|R\|$  denota o comprimento, área ou volume de  $R$ . Então  $R$  é ótima se, e somente se,  $\pi(\theta_1 \mid \mathbf{x}) \geq \pi(\theta_2 \mid \mathbf{x})$  para todo  $\theta_1 \in R$  e  $\theta_2 \notin R$ , salvo para conjuntos de medida nula.

# Estimação por Regiões e Teoria de decisão

**Prova:** Para todo  $R \in \mathcal{A}$  temos

$$\begin{aligned} E(L(R, \theta) \mid \mathbf{x}) &= \int_{\Theta} K\|R\| I(\theta \in R) \pi(\theta \mid \mathbf{x}) d\theta \\ &= K\|R\| \int_{\Theta} \pi(\theta \mid \mathbf{x}) I(\theta \in R) d\theta \\ &= K\|R\| \int_R \pi(\theta \mid \mathbf{x}) d\theta \\ &= K\|R\| \mathbb{P}(\theta \in R \mid \mathbf{x}) \\ &= K\|R\| \gamma \end{aligned}$$

$\Leftarrow$ ] Suponha que  $\pi(\theta_1 \mid \mathbf{x}) \geq \pi(\theta_2 \mid \mathbf{x})$  para todo  $\theta_1 \in R$  e  $\theta_2 \in R^c$

$\Rightarrow R$  está na região de maior massa probabilística *a posteriori*

$\Rightarrow$  Como para todo  $R \in \mathcal{A}$  temos que  $\mathbb{P}(\theta \in R \mid \mathbf{x})$ , então  $\|R\|$  é mínima entre todas as regiões  $R \in \mathcal{A}$

$\Rightarrow E(L(R, \theta) \mid \mathbf{x}) = K\|R\|\gamma$  é mínima  $\Rightarrow R$  é a região ótima.

# Estimação por Regiões e Teoria de decisão

$\Rightarrow$  Suponha que  $R$  é uma região ótima  $\Rightarrow E(L(R, \theta) \mid \mathbf{x})$  é mínima  $\Rightarrow K\|R\|^\gamma$  é mínima  $\Rightarrow \|R\|$  é mínima.

Como o volume sob  $\pi(\theta \mid \mathbf{x})$  restrito a  $R$  é fixo e igual a  $\gamma$  então  $\|R\|$  é mínimo se  $R$  contém todos os valores  $\theta_1 \in \Theta$  que têm mais alta densidade *a posteriori*. Além disto, todos os pontos  $\theta_2 \in R^c$  têm menos densidade *a posteriori* que qualquer um dos pontos  $\theta_1$  pertencentes à  $R$ .

## Regiões HPD: Algoritmo

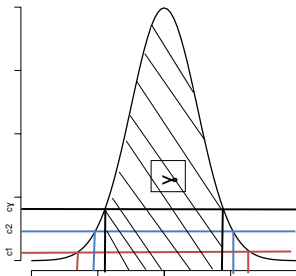
Na prática, não é muito simples obter as regiões HPD de forma exata. A construção do algoritmo passa por

**P1** : Construir uma rotina para determinar soluções para

$$\pi(\theta | x) = C,$$

para todo  $C$ , definindo as regiões

$$R_{\theta}^C(x) = \{\theta \in \Theta : \pi(\theta | x) \geq C\}$$



# Regiões HPD: Algoritmo

P2 : Construir uma rotina para calcular

$$P(\theta \in R_{\theta}^C(x) \mid x) = \int_{R_{\theta}^C(x)} \pi(\theta \mid x) d\theta.$$

P3 : A região HPD é determinada quando encontrarmos o maior  $C = C_{\gamma}$  tal que

$$P(\theta \in R_{\theta}^{C_{\gamma}}(x) \mid x) = \int_{R_{\theta}^{C_{\gamma}}(x)} \pi(\theta \mid x) d\theta \geq \gamma$$



## Regiões HPD: Algoritmo

Se a região HPD é um intervalo e se temos uma amostra  $\theta_1, \dots, \theta_n$  da distribuição *a posteriori*  $\pi(\theta | x)$ , então se pode aproximar os intervalos HPD como segue:

**Algoritmo (Chen and Shao (J. Comp Graph. Stat., 1999)):**

**P1 :** Ordene a amostra  $\theta_{(1)}, \dots, \theta_{(n)}$

**P2 :** Determine os intervalos com credibilidade  $\gamma$  fazendo

$$R_i(x) = [\theta_{(i)}; \theta_{i+\lfloor n\gamma \rfloor - 1}], \forall i = 1, \dots, n - \lfloor n\gamma \rfloor + 1.$$

**P3 :** O intervalo HPD é o intervalo  $R_{i_0}(x)$  tal que

$$\theta_{i_0+\lfloor n\gamma \rfloor - 1} - \theta_{i_0} = \min_i \{ \theta_{i+\lfloor n\gamma \rfloor - 1} - \theta_i, i = 1, \dots, n - \lfloor n\gamma \rfloor + 1 \}.$$

onde  $\lfloor a \rfloor$  é o maior inteiro menor ou igual que  $a$ .

## Regiões HPD: Algoritmo

**Exemplo:** Suponha que queiramos obter um intervalo com credibilidade  $\gamma = 0,8$

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$	$\theta_5$	$\theta_6$	$\theta_7$	$\theta_8$	$\theta_9$	$\theta_{10}$
1	3	8	8	9	12	14	18	26	30
$\theta_{(1)}$	$\theta_{(2)}$	$\theta_{(3)}$	$\theta_{(4)}$	$\theta_{(5)}$	$\theta_{(6)}$	$\theta_{(7)}$	$\theta_{(8)}$	$\theta_{(9)}$	$\theta_{(10)}$

Temos que  $n - \lfloor n\gamma \rfloor + 1 = 10 - \lfloor 10 \times 0.8 \rfloor + 1 = 3$

Candidatos

$$R_1 = [\theta_{(1)}, \theta_{(1+8-1)}] = [\theta_{(1)}, \theta_{(8)}] = [1, 18] \rightarrow \text{Amp} = 17, p = 0,8$$

$$R_2 = [\theta_{(2)}, \theta_{(2+8-1)}] = [\theta_{(2)}, \theta_{(9)}] = [3, 26] \rightarrow \text{Amp} = 23, p = 0,8$$

$$R_3 = [\theta_{(3)}, \theta_{(3+8-1)}] = [\theta_{(3)}, \theta_{(10)}] = [8, 30] \rightarrow \text{Amp} = 22, p = 0,8$$

O Intervalo HPD com credibilidade 0,8 é

$$R_1 = [1, 18]$$