

# AGA 0505 - Análise de Dados em Astronomia

## 7. Inferência Bayesiana

Laerte Sodr  Jr.

1o. semestre, 2024



# o teorema de Bayes

- o teorema de Bayes:
  - Bayes ~1740, Laplace 1774
  - $a$  e  $b$ : variáveis aleatórias
  - probabilidade conjunta de  $a$  e  $b$ , dado um certo modelo  $M$ :

$$\begin{aligned}P(a, b|M) &= P(a|b, M) \times P(b|M) = \\ &= P(b|a, M) \times P(a|M)\end{aligned}$$

e, portanto:

$$P(a|b, M) = \frac{P(b|a, M) \times P(a|M)}{P(b|M)}$$



# o teorema de Bayes e a análise de dados

## • a) inferência dos parâmetros de um modelo:

- temos um conjunto de dados  $D$  que queremos explicar com um modelo  $M$  que tem um conjunto de parâmetros  $w$
- objetivo: determinar  $w$
- exemplo:  $D = \{x_1, y_1, x_2, y_2, \dots, x_n, y_n\}$ ,  
 $M : y = a + bx + \epsilon$ ,  $w = \{a, b\}$
- teorema de Bayes ( $a \rightarrow w, b \rightarrow D$ ):

$$P(w|D, M) = \frac{P(D|w, M)P(w|M)}{P(D|M)}$$

$$P(D|M) = \int P(D|w, M)P(w|M)dw$$

## • b) comparação de modelos:

- temos dois modelos,  $M_1$  e  $M_2$ , para os dados  $D$ ; qual modelo escolher?
- exemplo:  $M_1$  : reta,  $M_2$  : parábola
- teorema de Bayes ( $a \rightarrow M, b \rightarrow D$ ):

$$P(M|D) = \frac{P(D|M)P(M)}{P(D)}$$

- comparação de  $M_1$  e  $M_2$ :

$$\frac{P(M_1|D)}{P(M_2|D)} = \frac{P(D|M_1) P(M_1)}{P(D|M_2) P(M_2)}$$

escolhemos o modelo mais provável











## exemplo: quantos peixes tem no lago?

- exemplo proposto por Rasmus Bååth U. de Lund ([www.sumsar.net](http://www.sumsar.net))
- Vamos a um lago e pescamos 7 peixes. Marcamos cada um e os devolvemos ao lago. Alguns dias depois voltamos ao lago, pescamos 20 peixes (dia de sorte!) e verificamos que 3 deles estão marcados. Quantos peixes tem no lago?
- Seja
  - T: número total de peixes no lago: ?
  - M: número de peixes marcados: 7
  - K: número de peixes pescados: 20
  - X: número de peixes pescados que estão marcados: 3

- Solução de MV: supomos que a fração de peixes marcados ( $X/K$ ) resultante de nossa medida (pescaria) é representativa
- regra de 3: se com  $K=20$  peixes temos  $X=3$  marcados, com T teremos  $M=7$ , logo:

$$\hat{T} = MK/X = 7 \times 20/3 \simeq 46.67$$



Lagoa Feia, Coração de Jesus, MG







## que prior escolher?

- em muitos casos usam-se *priors não-informativos*
- *priors não-informativos* para parâmetros de localização (ex.: média de uma gaussiana):

$$P(w) = \begin{cases} \frac{1}{b-a} & \text{se } a < w < b \\ 0 & \text{se } w < a \text{ ou } w > b \end{cases}$$

*prior uniforme* dentro de um intervalo de interesse

- nesse caso, o valor mais provável do posterior (MAP, *maximum a posteriori*) com prior uniforme é igual à solução de máxima verossimilhança
- outra opção: gaussiana muito larga

- *priors não-informativos* para parâmetros de escala (ex.: dispersão de uma gaussiana):

$$P(w) = \begin{cases} \frac{1}{w} & \text{se } a < w < b \\ 0 & \text{se } w < a \text{ ou } w > b \end{cases}$$

*prior de Jeffreys*: uniforme em  $\log w$

- os parâmetros que aparecem nos *priors* (como  $a$  e  $b$  acima) são chamados *hiperparâmetros*, para distinguir dos parâmetros do modelo
- *priors* podem ser próprios (normalizados) ou impróprios (não normalizados)





## exemplo 2: inferência da média de uma gaussiana com um prior gaussiano

- temos uma única observação,  $x_1$ , que queremos modelar como uma gaussiana de média  $\mu$  desconhecida e desvio padrão  $\sigma$  conhecido: queremos determinar  $\mu$
- a verossimilhança é  $N(\mu, \sigma)$
- vamos supor para  $\mu$  um prior gaussiano de média  $\mu_0$  e desvio padrão  $\tau_0$ ,  $P(\mu) \sim N(\mu_0, \tau_0)$
- posterior de  $\mu$ :  $P(\mu|x_1, \sigma) \propto P(x_1|\mu, \sigma)P(\mu) \propto \exp\left[-\frac{(x_1-\mu)^2}{2\sigma^2} - \frac{(\mu-\mu_0)^2}{2\tau_0^2}\right] \propto \exp\left[-\frac{(\mu-\mu_1)^2}{2\tau_1^2}\right] \sim N(\mu_1, \tau_1)$
- o posterior também é gaussiano com

$$\mu_1 = \frac{\frac{\mu_0}{\tau_0^2} + \frac{x_1}{\sigma^2}}{\frac{1}{\tau_0^2} + \frac{1}{\sigma^2}} \qquad \frac{1}{\tau_1^2} = \frac{1}{\tau_0^2} + \frac{1}{\sigma^2}$$

note que se  $\tau_0 \gg \sigma$ , a influência do prior é desprezível  
se  $\tau_0 \ll \sigma$ , o prior é muito *informativo*



## exemplo 2: inferência da média de uma gaussiana com um prior gaussiano

- agora temos  $n$  observações:  $D = \{x_i, \sigma\}$
- modelo: gaussiana de média  $\mu$  desconhecida e desvio padrão  $\sigma$  conhecido:  $P(x_i, \sigma | \mu) \sim N(\mu, \sigma)$
- prior de  $\mu$ :  $P(\mu) \sim N(\mu_0, \tau_0)$
- verossimilhança dos dados (supostos independentes):

$$P(D|\mu) = \prod_{i=1}^n P(x_i, \sigma|\mu) \propto \prod_{i=1}^n \exp\left[-\frac{1}{2\sigma^2}(x_i - \mu)^2\right] \propto \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right]$$

- posterior de  $\mu$ :  $P(\mu|D, \sigma) \sim N(\mu_n, \tau_n)$ , onde  $\mu_n = \frac{\frac{\mu_0}{\tau_0^2} + \frac{n\bar{x}}{\sigma^2}}{\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}}$ ,  $\frac{1}{\tau_n^2} = \frac{1}{\tau_0^2} + \frac{n}{\sigma^2}$  e  $\bar{x} = \frac{1}{n} \sum_i^n x_i$
- para  $n$  grande o efeito do prior fica desprezível!
- *o prior não deve dominar o resultado, a menos que seja relevante!*

## exemplo 3: regressão linear bayesiana

- temos um conjunto de dados  
 $D = \{x_i, y_i\}$ ,  $i = 1, \dots, n$  independentes e com a mesma variância ( $\sigma^2$ )
- modelo:  $y_i = a + bx_i + \epsilon_i$ ,  
com erros gaussianos:  $\epsilon_i \sim N(0, \sigma^2)$
- parâmetros:  $w = \{a, b, \sigma\}$
- objetivo: determinar as distribuições de probabilidades de  $a, b, \sigma$
- teorema de Bayes:  
 $P(a, b, \sigma | D) \propto P(D | a, b, \sigma) P(a, b, \sigma)$
- verossimilhança de uma medida:  
 $P(D_i | a, b, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{(y_i - a - bx_i)^2}{2\sigma^2} \right]$

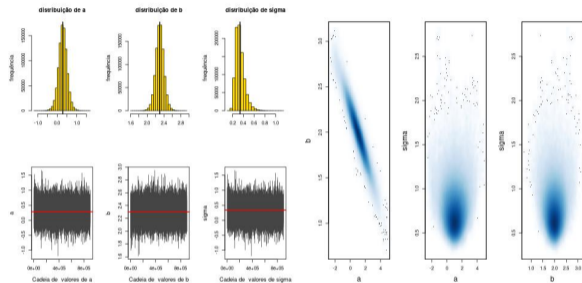
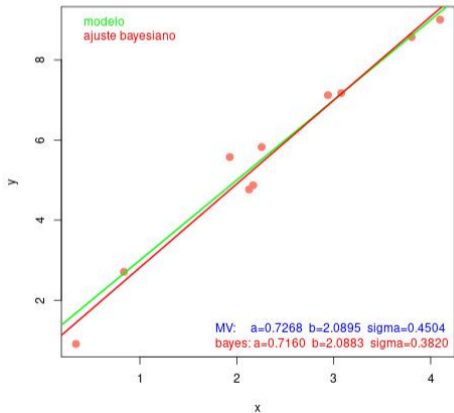
- verossimilhança da amostra:  
 $P(D | a, b, \sigma) = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp \left[ -\sum_i \frac{(y_i - a - bx_i)^2}{2\sigma^2} \right]$
- vamos supor priores uniformes para  $a$  e  $b$  e o prior de Jeffreys para  $\sigma$ :

$$P(a, b, \sigma) \propto \frac{1}{\sigma^2}$$

- posterior:  
 $P(a, b, \sigma | D) \propto \sigma^{-(n+2)} \exp \left[ -\sum_i \frac{(y_i - a - bx_i)^2}{2\sigma^2} \right]$

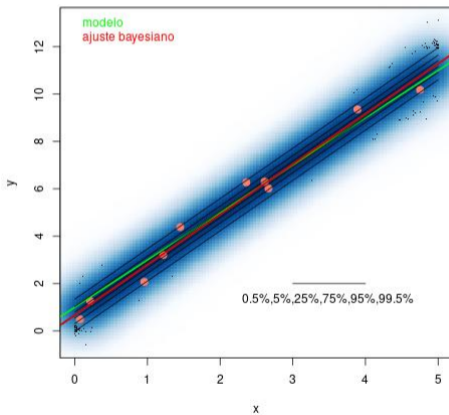
# exemplo 3: regressão linear bayesiana

M:  $y = 1 + 2x + \epsilon$



# exemplo 3: regressão linear bayesiana

**M:  $y = 1 + 2x + \text{eps}$**



## exemplo 4: o problema do farol

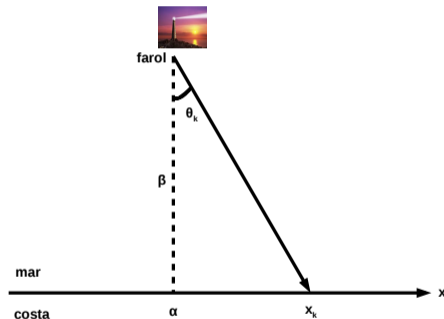
exemplo de uma modelagem bayesiana  
Sivia & Skilling, sec. 2.4 (Gull, 1988)

- um farol está em uma posição  $\alpha$  ao longo de uma costa reta e a uma distância  $\beta$  no mar;
- girando, ele emite uma série de pulsos curtos altamente colimados em intervalos de tempo aleatórios (e portanto em azimutes  $\theta$  também aleatórios);
- $N$  destes pulsos são detectados por sensores na costa, mas só as posições  $D = \{x_1, x_2, \dots, x_N\}$ , não as direções:

onde está o farol?

- objetivo: determinar  $\alpha$  e  $\beta$

$$P(\alpha, \beta | D) \propto P(D | \alpha, \beta) P(\alpha, \beta)$$



# exemplo 4: o problema do farol

- da geometria do problema:  $\text{tg}\theta = (x - \alpha)/\beta$
- vamos atribuir um prior uniforme para  $\theta$ :
  - para o pulso ser detectado, ele deve ter sido emitido com  $\theta$  em  $-\frac{\pi}{2} \leq \theta_k \leq \frac{\pi}{2}$ :  

$$P(\theta_k) = \frac{1}{\pi}$$
- mudança de variáveis:  $|P(x)dx| = |P(\theta)d\theta|$
- derivando em relação a  $x$ :

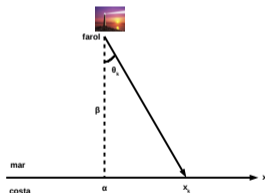
$$\frac{d \text{tg}\theta}{d\theta} = \sec^2 \theta = \frac{1}{\beta} \frac{dx}{d\theta},$$

logo,

$$\frac{d\theta}{dx} = \frac{1}{\beta \sec^2 \theta} = \frac{1}{\beta(1 + \text{tg}^2\theta)} = \frac{\beta}{\beta^2 + (x - \alpha)^2}$$

- como  $P(x|\alpha, \beta) = P(\theta|\alpha, \beta) \left| \frac{d\theta}{dx} \right|$  vem que  

$$P(x|\alpha, \beta) = \frac{\beta}{\pi[\beta^2 + (x - \alpha)^2]},$$
 que é uma distribuição de Cauchy!

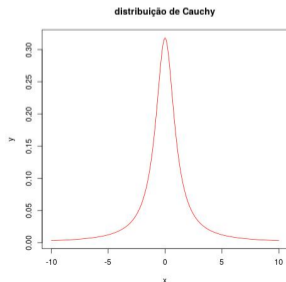


## exemplo 4: o problema do farol

- distribuição de Cauchy

$$P(x|\alpha, \beta) = \frac{\beta}{\pi[\beta^2 + (x - \alpha)^2]}$$

distribuição simétrica em relação a  $\alpha$  e com  $FWHM = 2\beta$  (FWHM: largura a meia altura)



- verossimilhança:

$$P(D|\alpha, \beta) = \prod_{k=1}^N P(x_k|\alpha, \beta)$$

e, portanto,

$$\ln P(D|\alpha, \beta) =$$

$$= \text{cte} + N \ln \beta - \sum_{k=1}^N \ln[\beta^2 + (x_k - \alpha)^2]$$

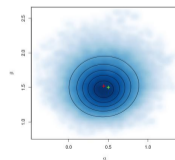
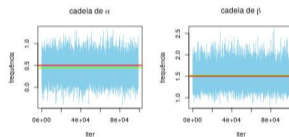
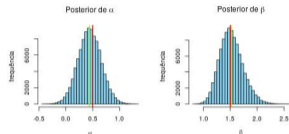
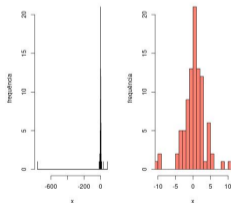
- vamos adotar priores não informativos para  $\alpha$  e  $\beta$ : nesse caso o posterior é proporcional à verossimilhança

# exemplo 4: o problema do farol

- log do posterior:

$$\ln P(\alpha, \beta | D) = \text{cte} + N \ln \beta - \sum_{k=1}^N \ln[\beta^2 + (x_k - \alpha)^2]$$

- simulação com  $N = 100$ ,  $\alpha = 0.5$ ,  $\beta = 1.5$



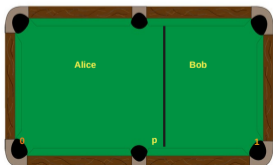
- inferência com MCMC:

$$\alpha = 0.44 \pm 0.21, \beta = 1.52 \pm 0.22$$



# quando estimativas bayesianas e frequentistas divergem?

- métodos bayesianos lidam com distribuições de probabilidades, métodos frequentistas lidam com estimativas de ponto
- em muitos casos o máximo do posterior (MAP) é idêntico à estimativa de MV
- em outros casos não; por exemplo:
  - quando se usa priores informativos
  - no trato de hiperparâmetros ou “nuisance parameters”

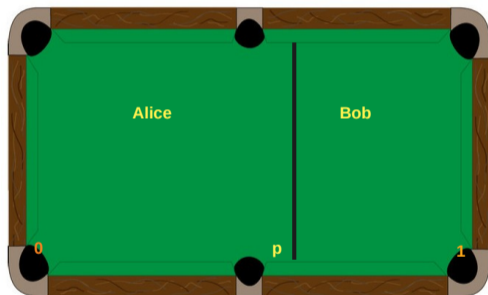


- exemplo: *um jogo de bilhar bayesiano* baseado no trabalho de 1763 de Bayes  
ver <http://jakevdp.github.io/blog/2014/06/06/frequentism-and-bayesianism-2-when-results-differ/>
- Carol joga bolas, de costas e sem viés, numa mesa de bilhar que tem uma marca: se elas caem de um lado da marca, Alice ganha um ponto, se caem do outro, Bob ganha um ponto; ganha o jogo quem primeiro fizer 6 pontos
- num certo jogo, após 8 bolas, Alice tem 5 pontos e Bob tem 3
- qual é a probabilidade de Bob ganhar o jogo?

# quando estimativas bayesianas e frequentistas divergem?

- após 8 bolas, Alice tem 5 pontos e Bob tem 3
- abordagem frequentista:
  - a probabilidade  $p$  da bola cair do lado da Alice é  $p = 5/8$
  - a probabilidade  $p$  da bola cair do lado de Bob é  $1 - p$
  - para Bob ganhar o jogo, ele tem que marcar 3 pontos seguidos
  - probabilidade disso:

$$P_{freq} = (1 - p)^3 = 0.0527$$



# quando estimativas bayesianas e frequentistas divergem?

- abordagem Bayesiana:
  - seja  $B$  o evento “Bob ganha”
  - dados:  $D = \{n_A, n_B\} = \{5, 3\}$
  - $p$ : probabilidade (desconhecida) que a bola caia na área de Alice
  - queremos  $P(B|D)$
  - note que o valor de  $p$  não interessa!  
ele é um *nuisance parameter*

- “sumimos” com  $p$  via marginalização:

$$\begin{aligned}P(B|D) &= \int P(B, p|D)dp = \\ &= \int P(B|p, D)P(p, D)dp = \\ &= \int P(B|p, D) \frac{P(D|p)P(p)}{P(D)} dp = \\ &= \frac{\int P(B|p, D)P(D|p)P(p)dp}{\int P(D|p)P(p)dp}\end{aligned}$$



# quando estimativas bayesianas e frequentistas divergem?

abordagem Bayesiana:

- marginalização:

$$P(B|D) = \frac{\int P(B|p, D)P(D|p)P(p)dp}{\int P(D|p)P(p)dp}$$

- para ganhar a partida, Bob tem que ganhar 3 jogadas seguidas:

$$P(B|p, D) = P(B|p) = (1 - p)^3$$

- vamos supor  $P(p)$  uniforme entre 0 e 1
- verossimilhança: distribuição binomial

$$P(D|p) \propto p^5(1 - p)^3$$

- logo,

$$P(B|D) = \frac{\int_0^1 (1 - p)^6 p^5 dp}{\int_0^1 (1 - p)^3 p^5 dp}$$

- a função beta é:

$$\beta(n, m) = \int_0^1 (1 - p)^{n-1} p^{m-1} dp$$

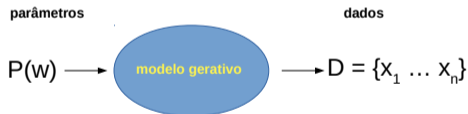
- logo,  $P(B|D) = \frac{\beta(6+1, 5+1)}{\beta(3+1, 5+1)} \simeq 0.091$
- compare com MV:  $P_{freq} = 0.053!$

# modelos gerativos

- modelos gerativos: inferência baseada em simulações dos dados
  - simulamos parâmetros a partir de um prior,  $P(w)$ ,
  - usamos esses parâmetros para simular “dados”, e
  - usamos os dados simulados para estimar o posterior dos parâmetros

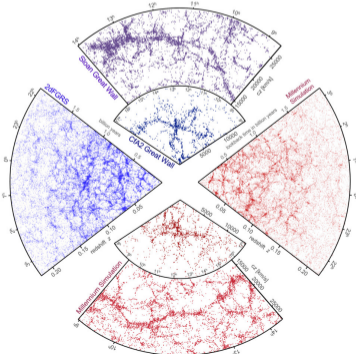
## modelos gerativos

Simulação de dados a partir do prior dos parâmetros e do modelo



# ABC: Approximate Bayesian Computation

- objetivo: estimar o posterior dos parâmetros de um modelo
  - a) sem calcular a verossimilhança
  - b) mas usando um processo gerativo de gerar dados a partir de parâmetros



- porquê?
  - há situações onde a verossimilhança é intratável
  - exemplo: estimativa de parâmetros cosmológicos comparando observações e simulações da distribuição de galáxias
- como proceder?
  - ao invés de se comparar os dados (e.g. posições de galáxias) diretamente, pode-se comparar estatísticas que “resumem” propriedades importantes dos dados tanto nas observações quanto nas simulações
- exemplos: distância média entre galáxias, variância do número de galáxias em esferas de raio 8 Mpc, ...

# ABC: Approximate Bayesian Computation

- objetivo: estimar o posterior dos parâmetros  $w$  de um modelo, usando um processo gerativo a partir dos priores  $P(w)$
- ABC: apenas dados simulados que concordam com as observações dentro de uma certa “tolerância” são considerados para amostrar o posterior



- algoritmo ABC:
  - 1. amostre  $w_p$  do prior  $P(w)$
  - 2. simule dados  $D_p$  com  $w_p$
  - 3. calcule as estatísticas que resumizam os dados:  $x_p = \text{resumo}(w_p)$
  - 4. aceite  $w_p$  se  $|x_p - x_{obs}| < \epsilon$  (tolerância)
  - 5. retorne a 1
- atenção: em geral (mas nem sempre),
  - dados contínuos: tolerância  $\epsilon > 0$
  - dados discretos: tolerância  $\epsilon = 0$







