

## 2.15 Estimações de parâmetros

(89)

O problema das estimações de parâmetros pode ser definido como: dadas medições  $y_k$  relacionadas com variáveis discritas  $x_k$  tais que

$$y_k = f(x_k, \theta)$$

com  $f(\cdot, \theta)$  sendo uma função de parâmetros  $\theta$  (escalar ou vetor), determinar uma estimativa  $\hat{\theta}$  de  $\theta$  de modo que  $\hat{\theta}$  seja o mais próximo possível de  $\theta$ . As soluções a este problema serão classificadas em três abordagens distintas: quadrados mínimos, máximos de verossimilhança e estimações bayesiana.

Quanto ao problema, esse pode ser classificado em linear ou não-linear, dependendo da forma de  $f(\cdot, \theta)$  com relação a  $\theta$ .

Exemplo: estimações lineares ou não-lineares?

Se  $f(x_k, \theta)$  puder ser escrito na forma

$$f(x_k, \theta) = \varphi(x_k)^T \cdot \theta$$

com  $\varphi(\cdot)^T: \mathbb{R}^m \rightarrow \mathbb{R}^m$  sendo um vetor coluna que é função de  $x_k \in \mathbb{R}^m$  e  $\theta \in \mathbb{R}^m$  sendo o vetor de parâmetros, então o problema de estimação é de tipo linear. Se  $f(x_k, \theta) = a \cdot \sin(x_k) + b \cdot \cos(x_k)$ , então

$$\varphi(x_k)^T = [\sin(x_k) \quad \cos(x_k)]$$

$$\theta^T = [a \quad b]$$

Há se  $f(x_k, \theta) = a \cdot \sin(x_k) + b \cdot \cos(x_k) + \delta$ , que é uma função afim mas pode ser escrita como  $\varphi(x_k)^T \cdot \theta$ . Mas, com a seguinte mudança:

$$\begin{aligned} y_{k-1} &= a \cdot \sin(x_k) + b \cdot \cos(x_k) \\ &= \underbrace{[\sin(x_k) \quad \cos(x_k)]}_{\varphi(x_k)^T} \cdot \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{\theta} \end{aligned}$$

O problema se torna linear.

Com  $f(x_k) = \cos(x_k + \theta)$ , temos então um problema não-linear.

Algunas propriedades desejadas para estimadores são:

i) Estimativa não-tendenciosa:

$$E\{\hat{\theta}\} = \theta$$

ii) Estimativa não-tendenciosa da variância mínima:

Para toda estimativa  $\theta^*$  que satisfaça

$$E\{\theta^*\} = \theta$$

$\hat{\theta}$  é uma estimativa não-tendenciosa da variância mínima se  $E\{\hat{\theta}\} = \theta$  e

$$\text{Var}\{\hat{\theta}\} \leq \text{Var}\{\theta^*\}.$$

iii) Estimativa coerente

$$\lim_{K \rightarrow \infty} \Pr\{|\hat{\theta} - \theta| > \epsilon\} = 0$$

para todo  $\epsilon > 0$ . Assim, uma estimativa é

$$\text{coerente se } \lim_{n \rightarrow \infty} E\{\hat{\theta}\} = \theta \text{ e } \lim_{n \rightarrow \infty} \text{Var}\{\hat{\theta}\} = 0$$

As definições acima se aplicam a problemas de estimativas determinísticas, em que  $\theta$  é considerado uma constante (i.e.,  $\text{Var}\{\theta\} = 0$ ). Em estimativas estocásticas temos ainda

iv) Estimativa consistente

$$E\{\hat{\theta} - \theta\} = 0 \text{ e } \text{Var}\{\hat{\theta}\} = \text{Var}\{\theta\}$$

para uma estimativa inconsistente, basta  $E\{\hat{\theta} - \theta\} \neq 0$

v) Estimativa consistente pessimista

$$E\{\hat{\theta} - \theta\} = 0 \text{ e } \text{Var}\{\hat{\theta}\} > \text{Var}\{\theta\}$$

também chamada de estimativa conservativa.

vi) Estimativa consistente otimista

$$E\{\hat{\theta} - \theta\} = 0 \text{ e } \text{Var}\{\hat{\theta}\} < \text{Var}\{\theta\}$$

Em muitas aplicações, estimativas otimistas são perigosas quando são usadas em processo de decisão.

## 2.15.1 Estimativas por mínimos quadrados

Nos métodos de mínimos quadrados, o problema se escreve da forma seguinte:

$$\hat{\theta} = \arg \min_{\theta} J(r, \theta)$$

com  $r$  sendo os resíduos (ou erros) de estimativas,

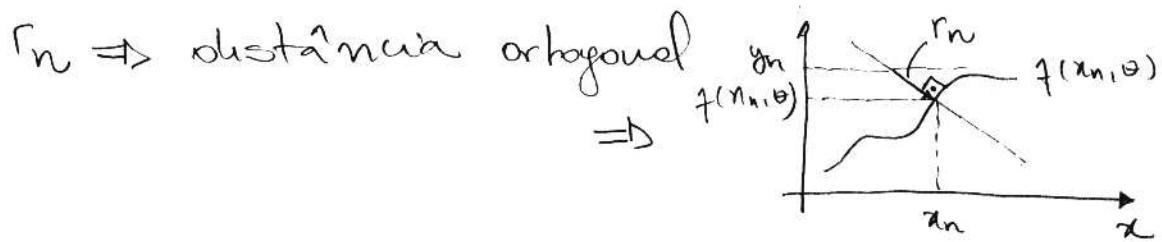
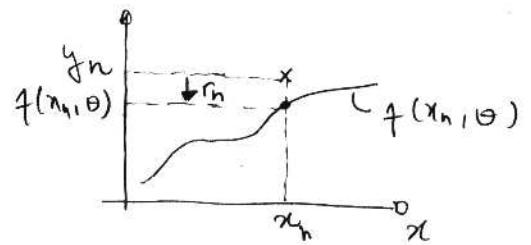
$J(r, \theta)$  é uma função do hpo

$$J(r, \theta) = \sum_{n=1}^N r_n(x_n, y_n, \theta)^2$$

Assim,  $\hat{\theta}$  é a estimativa que minimiza a soma dos quadrados dos resíduos. Em função da forma como os resíduos são definidos, temos os mínimos quadrados (MQ) ou ainda os mínimos quadrados totais (MQT). Em MQ, os resíduos são simplesmente

$$r_n = y_n - f(x_n, \theta) \Rightarrow$$

enquanto que, no MQT,



Exemplo: Estimações dos parâmetros de uma reta

Siga o modelo

$$y = a \cdot x + b$$

Este modelo pode ser escrito como

$$y_n = \varphi(x_n)^T \cdot \theta, \text{ com } \varphi(x_n)^T = [x_n \ 1]$$

e  $\theta = [a \ b]^T$ . Pela abordagem MQ, temos

$$r_n = y_n - a \cdot x_n - b$$

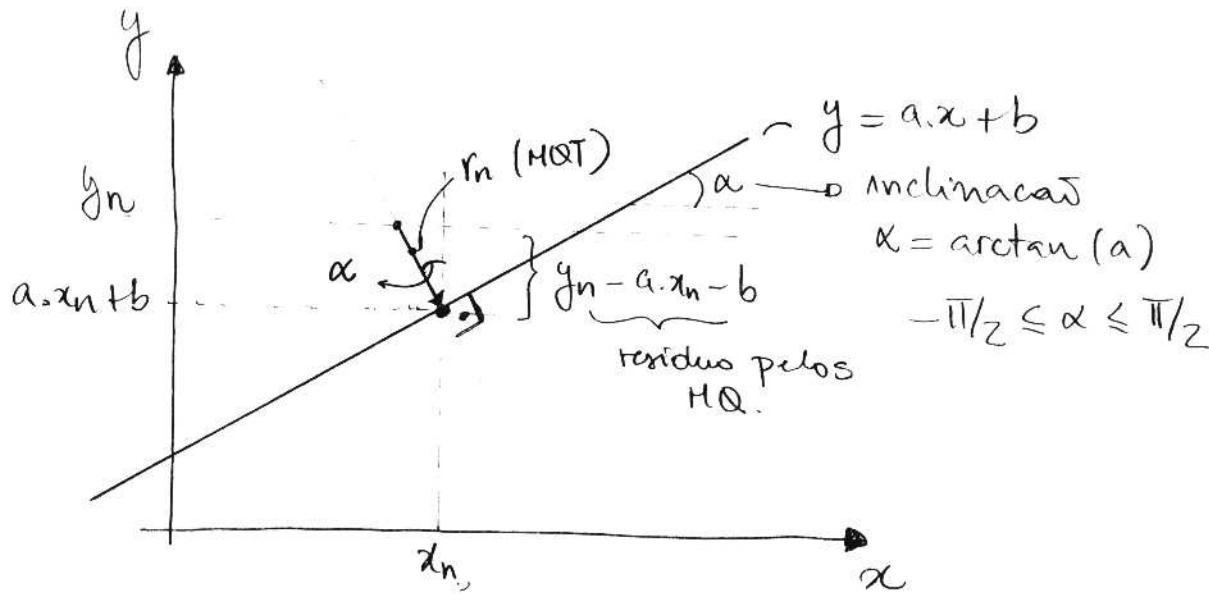
enquanto que, com a abordagem MQT, os resíduos

(43)

são definidos como

$$r_n = (y_n - (a \cdot x_n + b)) / \cos(\alpha)$$

com  $\alpha = \arctan(a)$ . O esboço de  $r_n$  é mostrado abaixo.



Pelo método dos mínimos quadrados, a solução deste problema é dada pelo mínimo de  $J(r, \theta)$  com relação a  $\theta$ , que satisfaz a

$$\frac{\partial J(r, \theta)}{\partial \theta} \Big|_{\theta=\hat{\theta}} = 0$$

$$J(r, \theta) = \sum_{n=1}^N (y_n - a \cdot x_n - b)^2 = \sum_{n=1}^N r_n^2$$

$$\frac{\partial J}{\partial \theta} = \begin{bmatrix} \frac{\partial J}{\partial a} \\ \frac{\partial J}{\partial b} \end{bmatrix} = \sum_{n=1}^N 2 \cdot \frac{\partial r_n}{\partial \theta} \cdot r_n$$

$$\frac{\partial J}{\partial a} = \sum_{n=1}^N 2 \cdot (y_n - a \cdot x_n - b) \cdot (-x_n)$$

$$\frac{\partial J}{\partial b} = \sum_{n=1}^N 2 \cdot (y_n - a \cdot x_n - b) \cdot (-1)$$

$$\frac{\partial J}{\partial a} \Big|_{\begin{array}{l} a = \hat{a} \\ b = \hat{b} \end{array}} = 0 \Rightarrow \sum_{n=1}^N (\hat{a} \cdot x_n^2 + \hat{b} \cdot x_n - x_n \cdot y_n) = 0$$

$$\Rightarrow \hat{a} \cdot \sum_{n=1}^N x_n^2 + \hat{b} \cdot \sum_{n=1}^N x_n = \sum_{n=1}^N x_n \cdot y_n$$

$$\Rightarrow \hat{a} \cdot \frac{1}{N} \cdot \sum_{n=1}^N x_n^2 + \hat{b} \cdot \frac{1}{N} \cdot \sum_{n=1}^N x_n = \frac{1}{N} \sum_{n=1}^N x_n \cdot y_n$$

$$\frac{\partial J}{\partial b} \Big|_{\begin{array}{l} a = \hat{a} \\ b = \hat{b} \end{array}} = 0 \Rightarrow \sum_{n=1}^N (\hat{a} \cdot x_n + \hat{b} - y_n) = 0$$

$$\Rightarrow \hat{a} \cdot \sum_{n=1}^N x_n + \hat{b} \cdot \sum_{n=1}^N 1 = \sum_{n=1}^N y_n$$

$$\Rightarrow \hat{a} \cdot \frac{1}{N} \cdot \sum_{n=1}^N x_n + \hat{b} = \frac{1}{N} \cdot \sum_{n=1}^N y_n$$

Definindo:

$$\bar{x} = \frac{1}{N} \cdot \sum_{n=1}^N x_n \quad ; \quad \bar{y} = \frac{1}{N} \cdot \sum_{n=1}^N y_n$$

$$C_{xx} = \frac{1}{N} \cdot \sum_{n=1}^N x_n^2 \quad ; \quad C_{xy} = \frac{1}{N} \cdot \sum_{n=1}^N x_n \cdot y_n$$

Temos que resolver o sistema abaixo para encontrar as estimativas  $\hat{a}$  e  $\hat{b}$ :

$$\hat{a} \cdot C_{xx} + \hat{b} \cdot \bar{x} = C_{xy}$$

$$\hat{a} \cdot \bar{x} + \hat{b} = \bar{y}$$

fazendo  $\hat{b} = \bar{y} - \hat{a} \cdot \bar{x}$ , temos

$$\hat{a} \cdot C_{xx} + \bar{x} \cdot \bar{y} - \hat{a} \cdot \bar{x}^2 = C_{xy}$$

e assim

$$\hat{a} = \frac{C_{xy} - \bar{x} \cdot \bar{y}}{C_{xx} - \bar{x}^2} = \frac{S_{xy}}{S_{xx}}$$

$$\hat{b} = \bar{y} - \frac{S_{xy}}{S_{xx}} \cdot \bar{x}$$

$$\text{com } S_{xy} = \frac{1}{N} \cdot \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y}) = C_{xy} - \bar{x}\bar{y}$$

$$\text{e } S_{xx} = \frac{1}{N} \cdot \sum_{n=1}^N (x_n - \bar{x})^2 = C_{xx} - \bar{x}^2$$

Pode-se perceber que se  $S_{xx} = 0$  e  $S_{xy} \neq 0$ ,  
 então  $\hat{a} \rightarrow \infty$  (reta vertical com medições perfeitas). Neste  
 caso, a parametrização  $y = ax + b$  é limitada.  
 Melhor seria usar o modelo

$$a \cdot y + b \cdot x + c = 0$$

com resíduo  $r_n = a \cdot y_n + b \cdot x_n + c$ .

No caso dos MQT, a solução de

$$\hat{\theta} = \arg \min_{\theta} \sum_{n=1}^N \{(y_n - (a \cdot x_n + b)) / \cos(\arctan(a))\}^2$$

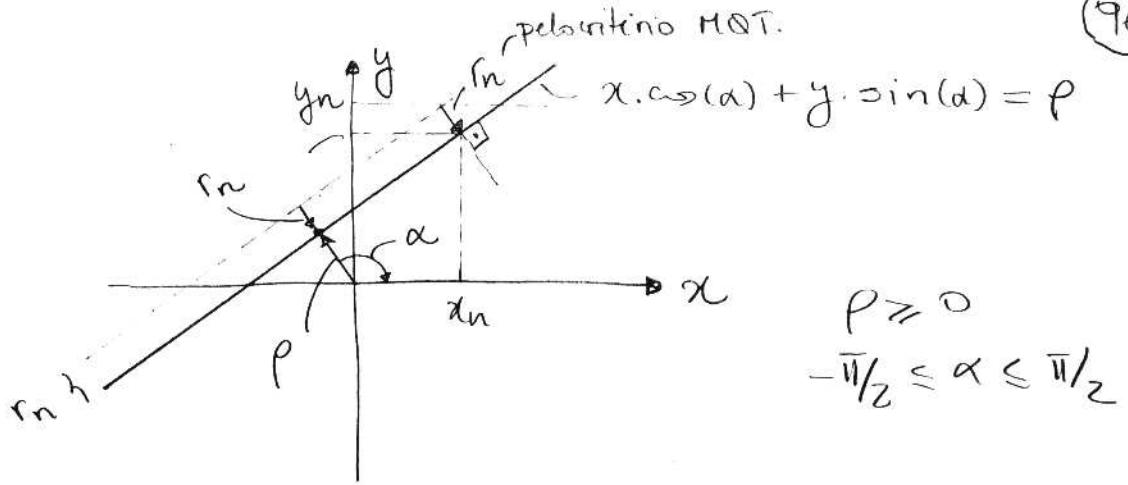
parece um pouco complicado, mas é resolvível. Assim,  
 assim, a escolha do modelo  $y = ax + b$  não é  
 adequada com retas verticais. O modelo

$$ay + bx + c = 0$$

é mais adequado e se aplica a qualquer tipo  
 de reta. No entanto, o resíduo pelo método MQT  
 parece complexo. Para reduzir a complexidade, o  
 modelo de reta dado por

$$\rho = x \cdot \cos(\alpha) + y \cdot \sin(\alpha)$$

é mais interessante pois envolve apenas dois  
 parâmetros ( $\rho$  e  $\alpha$ ).



Temos que  $\rho$  e  $\alpha$  são as coordenadas polares do ponto da reta que está mais próximo da origem

Pode-se verificar que, para o HQT, temos

$$\rho_n = \rho - x_n \cdot \cos(\alpha) - y_n \cdot \sin(\alpha)$$

E a minimização de

$$J = \sum_{n=1}^N (\rho - x_n \cdot \cos(\alpha) - y_n \cdot \sin(\alpha))^2$$

é obtida com parâmetros  $\hat{\rho}$ ,  $\hat{\alpha}$  tais que

$$\frac{\partial J}{\partial \rho} \Big|_{\begin{subarray}{l} \rho = \hat{\rho} \\ \alpha = \hat{\alpha} \end{subarray}} = \frac{\partial J}{\partial \alpha} \Big|_{\begin{subarray}{l} \rho = \hat{\rho} \\ \alpha = \hat{\alpha} \end{subarray}} = 0$$

Que resulta em:

$$\begin{aligned} \hat{\rho} &= \bar{x} \cdot \cos(\hat{\alpha}) + \bar{y} \cdot \sin(\hat{\alpha}) \\ \hat{\alpha} &= \frac{1}{2} \cdot \arctan \left( \frac{-2 \cdot S_{xy}}{S_{yy} - S_{xx}} \right) \end{aligned}$$

com  $\bar{x}$ ,  $\bar{y}$ ,  $S_{xy}$  e  $S_{xx}$  já definidos acima, e

$$S_{yy} = \frac{1}{N} \cdot \sum_{n=1}^N (y - \bar{y})^2$$

A generalização dos modelos apresentados a pouco pode ser dada pelo modelo linear abaixo:

$$y_n = \varphi(x_n)^T \cdot \theta + \epsilon_n$$

com  $y_k \in \mathbb{R}$  e  $\epsilon_k$  sendo o erro de medição. Então

$$J = \sum_{n=1}^N (y_n - \varphi(x_n)^T \cdot \theta)^2$$

Seu mínimo é obtido fazendo

$$\begin{aligned} \frac{\partial J}{\partial \theta} \Big|_{\theta=\hat{\theta}} &= 0 \\ &= \sum_{n=1}^N 2 \cdot (-\varphi(x_n)) \cdot (y_n - \varphi(x_n)^T \cdot \hat{\theta}) \end{aligned}$$

Então

$$\sum_{n=1}^N y_n \cdot \varphi(x_n) = \sum_{n=1}^N \varphi(x_n) \cdot \varphi(x_n)^T \cdot \hat{\theta}$$

Assim, a estimativa  $\hat{\theta}$  é dada por:

$$\left( \sum_{n=1}^N \varphi(x_n) \cdot \varphi(x_n)^T \right) \cdot \hat{\theta} = \sum_{n=1}^N \varphi(x_n) \cdot y_n$$

$m \times m$  simétrica

$m = \dim \hat{\theta}$

$1 \times m$  escalar

$$\hat{\theta} = \underbrace{\left( \sum_{n=1}^N \varphi(x_n) \cdot \varphi(x_n)^T \right)^{-1}}_{m \times m} \cdot \underbrace{\left( \sum_{n=1}^N \varphi(x_n) \cdot y_n \right)}_{m \times 1}$$

(98)

Exemplo: Aproximações polinomial

Dados pares de medições  $\{(x_i, y_i)\}$ , encontrar os parâmetros  $\theta = [a_0 \ a_1 \ \dots \ a_m]^T$  para o seguinte modelo polinomial

$$y = a_0 + a_1 \cdot x + a_2 \cdot x^2 + \dots + a_m \cdot x^m \\ = \varphi^T(x) \cdot \theta$$

com  $\varphi^T(x) = [1 \ x \ x^2 \ \dots \ x^m]$ .

Exemplo: Identificação de sistemas

Considere o seguinte modelo ARX:

$$y(k) = b_1 \cdot u(k-1) + \dots + b_m \cdot u(k-m) \\ - a_1 \cdot y(k-1) - \dots - a_m \cdot y(k-1)$$

Este modelo relaciona amostras da entrada e saída de um sistema cuja função de transferência discuta é dada por

$$\frac{y(k)}{u(k)} = \frac{b_1 z^{-1} + \dots + b_m z^{-m}}{1 + a_1 z^{-1} + \dots + a_m z^{-m}}$$

Considerando  $\theta = [b_1, \dots, b_m, a_1, \dots, a_m]^T$  e  $\varphi^T(u(k)) = [u(k-1) \ \dots \ u(k-m) \ -y(k-1) \ \dots \ -y(k-m)]$ ,  $\hat{\theta}$  pode ser obtido pelo método dos mínimos quadrados.

Observar que, quando

$$\epsilon_n = y_n - \varphi(x_n) \cdot \theta$$

se considerarmos que  $\epsilon_n \sim N(0, \sigma_\epsilon^2)$ , então

$$\sigma_\epsilon^2 = \frac{1}{N-1} \cdot \sum_{n=1}^N \epsilon_n^2 = \frac{1}{N-1} \cdot J(\theta)$$

valores de funções  
de ousto p/  $\hat{\theta}$ .

A estimativa  $\hat{\theta}$  pode ser escrita

$$\hat{\theta} = (\phi^\top \cdot \phi)^{-1} \cdot \phi^\top \cdot y$$

com  $\phi = \begin{bmatrix} \varphi(x_1)^\top \\ \vdots \\ \varphi(x_N)^\top \end{bmatrix}$  e  $y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$

$\underbrace{\qquad\qquad\qquad}_{N \times m}$

é o modelo de erro

$$Y = \phi \cdot \theta + \epsilon$$

com  $\epsilon = [\epsilon_1 \dots \epsilon_N]^\top$ ,  $\theta$  sendo o verdadeiro vetor de parâmetros e  $\epsilon \sim N(0, P_\epsilon)$  com

$$P_\epsilon = E\{\epsilon \cdot \epsilon^\top\} = \sigma_\epsilon^2 \cdot \mathbb{I}_N$$

considerando  $\epsilon_i$  e  $\epsilon_j$ ,  $i \neq j$ , independentes, e  $\mathbb{I}_N$  a matriz identidade de dimensão  $N$ . Então  $\epsilon_n$  é iid. O estimador  $\hat{\theta} = (\phi^\top \cdot \phi)^{-1} \cdot \phi^\top \cdot y$  tem as seguintes propriedades:

$$\begin{aligned} E\{\hat{\theta} - \theta\} &= E\{(\phi^\top \phi)^{-1} \cdot \phi^\top \cdot (\underbrace{\phi \cdot \theta + \epsilon}_{Y}) - \theta\} \\ &= E\{(\phi^\top \phi)^{-1} \cdot \phi^\top \cdot \phi \cdot \theta + (\phi^\top \phi)^{-1} \cdot \phi^\top \cdot \epsilon - \theta\} \\ &\quad \text{Im pois } A \cdot A^{-1} = \mathbb{I} \end{aligned}$$

$$= E\{\theta - \hat{\theta} + (\phi^\top \phi)^{-1} \cdot \phi^\top \varepsilon\}$$

$$= (\phi^\top \phi)^{-1} \cdot \phi^\top E\{\varepsilon\} = 0$$

ou seja,  $\hat{\theta}$  é não tendencioso (observe que  $\theta$  é tratado como uma constante). Assim,

$$E\{\hat{\theta} - \theta\} = E\{\hat{\theta}\} - E\{\theta\} = E\{\hat{\theta}\} - \theta = 0$$

significando que  $E\{\hat{\theta}\} = \theta$ . Então

$$\begin{aligned} P_{\hat{\theta}} &= E\{(\hat{\theta} - E\{\hat{\theta}\}).(\hat{\theta} - E\{\hat{\theta}\})^\top\} \\ &= E\{(\hat{\theta} - \theta).(\hat{\theta} - \theta)^\top\}, \text{ sendo } \hat{\theta} - \theta = (\phi^\top \phi)^{-1} \cdot \phi^\top \varepsilon \\ &= E\{(\phi^\top \phi)^{-1} \cdot \phi^\top \varepsilon \cdot \varepsilon^\top \cdot \phi \cdot (\phi^\top \phi)^{-1}\} \\ &= (\phi^\top \phi)^{-1} \cdot \phi^\top \underbrace{P_\varepsilon}_{\sigma_\varepsilon^2 \cdot \mathbb{I}_N} \cdot \phi \cdot (\phi^\top \phi)^{-1} \\ &= \underbrace{(\phi^\top \phi)^{-1} \cdot \phi^\top \phi}_{\mathbb{I}_m} \cdot (\phi^\top \phi)^{-1} \cdot \sigma_\varepsilon^2 \end{aligned}$$

$$P_{\hat{\theta}} = (\phi^\top \phi)^{-1} \cdot \sigma_\varepsilon^2 = \left( \sum_{n=1}^N \psi(x_n) \psi^\top(x_n) \right)^{-1} \cdot \sigma_\varepsilon^2$$

com  $\sigma_\varepsilon^2$  podendo ser estimado a partir de  $J(\hat{\theta})$

$$\sigma_\varepsilon^2 = \frac{1}{N-1} \cdot J(\hat{\theta}) \quad \left. \begin{array}{l} \text{com } J(\hat{\theta}) \leq J(\theta^*) \text{ pois } \hat{\theta} \text{ minimiza} \\ J, \text{ então } P_{\hat{\theta}} \leq P_{\theta^*}, \text{ e } \hat{\theta} \text{ é também} \\ \text{conhecido como estimador de mínima} \\ \text{variancia.} \end{array} \right.$$

pois  $\hat{\theta}$  é uma estimativa não tendenciosa de  $\theta$ .

Conforme foi exposto, o método dos mínimos quadrados considera que todas as medições têm pesos idênticos (i.e., mesma importância). No entanto

(10)

um  $y_i$  mais confiável do que um  $y_j$  pode ter um peso  $w_i > w_j$ . Isto resulta nos mínimos quadrados ponderados (MQP) :

$$J = \sum_{n=1}^N w_n \cdot (y_n - \varphi^T(x_n) \cdot \theta)^2 \\ = (\gamma - \Phi \cdot \theta)^T \cdot W \cdot (\gamma - \Phi \cdot \theta)$$

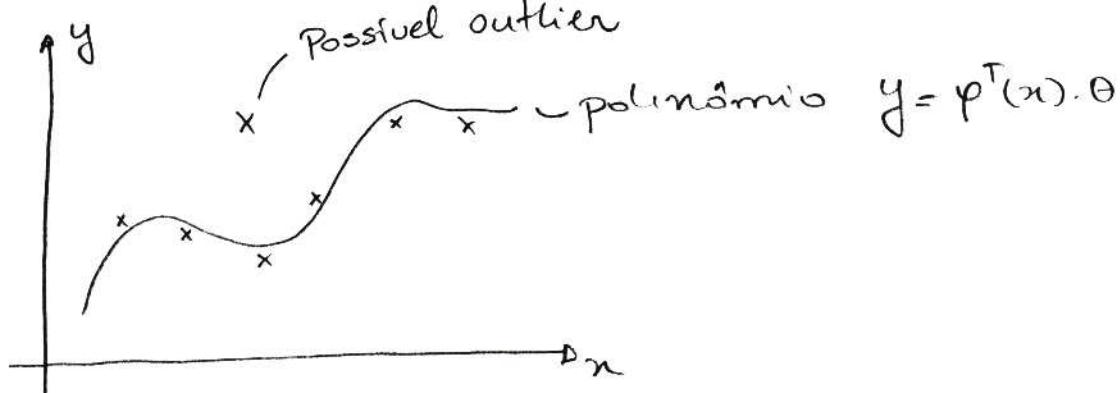
com  $W = \text{diag}(w_1, \dots, w_N)$ , então podemos verificar que

$$\frac{\partial J}{\partial \theta} \Big|_{\theta=\hat{\theta}} = \sum_{n=1}^N -w_n \cdot (y_n - \varphi^T(x_n) \cdot \hat{\theta}) \cdot \varphi^T(x_n) = 0$$

tem como soluções

$$\hat{\theta} = \left( \sum_{n=1}^N \varphi(x_n) \cdot \varphi(x_n)^T \cdot w_n \right)^{-1} \cdot \left( \sum_{n=1}^N \varphi(x_n) \cdot w_n \cdot y_n \right) \\ = (\Phi^T \cdot W \cdot \Phi)^{-1} \cdot \Phi^T \cdot W \cdot \gamma$$

Se  $w_n$  for escolhido apropriadamente, então a  $n$ -ésima medição pode ser descartada de estimações se ela for um "outlier", ou seja, uma medição incompatível com o modelo:



Podemos perceber que o mínimos quadrados é um caso particular do MOP, com  $w_n=1$ , ou seja não fazendo distinções entre bons e máis mediócos. Entretanto, o MQ é dito um estimador não-robusto pois basta um outlier para polarizar a estimativa. Uma abordagem robusta consiste em usar o seguinte algoritmo:

- i) Iniciar todos os pesos  $w_n$  iguais a 1 (MQ)
- ii) Lountains de iterações  $K=0$ .
- iii) Calcular  $\hat{\theta}_K$  usando MOP
- iv) Re-calcular  $w_n = \mu(r_n)$
- v) Incrementar  $K$  e retornar os passo (iii) se o critério de parada não for verificado

O algoritmo auma coluna pesos  $w_n$  de forma iterativa usando uma funcional  $\mu(r_n)$ , com  $r_n$  sendo o resíduo de estimativas usando a estimativa da  $K$ -ésima iteração:

$$r_n = y_n - \Phi(x_n)\hat{\theta}_K$$

Uma destas funcionais é a função de Huber

$$w_n = \min \left( 1, \frac{c \cdot \text{MAD}}{\text{abs}(r_n)} \right)$$

com  $c = 1,345$  e  $\text{MAD} = \text{MED}(|r_n - \text{MED}(r_n)|)$ . Os estimados obtidos com a função de Huber pertence à classe dos M-estimadores robustos.

## 2.15.2. Máximo de Verosimilhança

Dadas medições

$$y = f(x, \theta), \quad (\text{MV})$$

estimador do máximo de verosimilhança é dado por

$$\hat{\theta} = \arg \max_{\theta} P(y|\theta)$$

com  $P(y|\theta)$  sendo a função de dens. condicional de  $y$  dado  $\theta$ , também chamada de função de verosimilhança (likelihood function).

Maximizar  $P(y|\theta)$  é equivalente a maximizar  $\ln(P(y|\theta))$ , uma vez que  $\ln(\cdot)$  é uma função monótona crescente para  $P(y|\theta) > 0$ :

$$\hat{\theta} = \arg \max_{\theta} \ln(P(y|\theta))$$

Esta versão do MV é interessante quando  $P(y|\theta)$  envolve exponentiais.

Exemplo: Seja  $y = \varphi^T(x) \cdot \theta$  um modelo linear em  $\theta$  do qual pares  $\{y_n, x_n\}$  são obtidos (estimação linear) se as medições  $y_n$  são consideradas iid  $\sim N\{E[y_n], \sigma_y^2\}$  com  $E[y_n] = \varphi^T(x_n) \cdot \theta$ , uma vez que  $\theta$  não é uma var, mas sim uma constante. Então

$$P(y_n|\theta) = \frac{1}{\sqrt{2\pi} \cdot \sigma_y} \cdot \exp \left\{ -\frac{1}{2} \frac{(y_n - \varphi^T(x_n) \cdot \theta)^2}{\sigma_y^2} \right\}$$

Sendo  $p(y|\theta) = p(y_1, y_2, \dots, y_N | \theta)$  e  $y_n$  iid,  
então

$$\begin{aligned} p(y|\theta) &= \prod_{n=1}^N p(y_n|\theta) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma_y} \cdot \exp\left(-\frac{1}{2} \cdot \frac{(y_n - \varphi(x_n)^T \cdot \theta)^2}{\sigma_y^2}\right) \\ &= \frac{1}{(\sqrt{2\pi}\sigma_y)^N} \cdot \exp\left(\sum_{n=1}^N -\frac{1}{2} \frac{(y_n - \varphi(x_n)^T \cdot \theta)^2}{\sigma_y^2}\right) \end{aligned}$$

Como  $p(y|\theta)$  é uma função exponencial, é melhor maximizar  $\ln(p(y|\theta))$ :

$$\ln(p(y|\theta)) = -N \underbrace{\ln(\sqrt{2\pi}\sigma_y)}_{\text{Termo constante}} - \frac{1}{2\sigma_y^2} \sum_{n=1}^N (y_n - \varphi(x_n)^T \cdot \theta)^2$$

único termo que depende de  $\theta$

O máximo de  $\ln(p(y|\theta))$  é obtido fazendo a

$$\begin{aligned} \frac{\partial \ln(p(y|\theta))}{\partial \theta} \Big|_{\theta=\hat{\theta}} &= \frac{1}{\sigma_y^2} \cdot \sum_{n=1}^N (y_n - \varphi(x_n)^T \cdot \hat{\theta}) \cdot \varphi(x_n)^T \\ &= 0 \end{aligned}$$

Que significa, deve-se resolver

$$\sum_{n=1}^N (y_n - \varphi(x_n)^T \cdot \hat{\theta}) \cdot \varphi(x_n)^T = 0$$

resultando em

$$\hat{\theta} = \left( \sum_{n=1}^N \varphi(x_n) \cdot \varphi(x_n)^T \right)^{-1} \cdot \left( \sum_{n=1}^N \varphi(x_n)^T \cdot y_n \right)$$

que é o mesmo resultado obtido pelos mínimos quadrados.

(102)

Isto significa que, no caso de estimativas linear com erro de medições gaussiano de média nula e iid, MV e ML possuem resultados idênticos

Exemplo: Considere o seguinte modelo não-linear

$$y = \theta^2 + \epsilon$$

com  $\theta \in \mathbb{R}$  sendo uma constante a ser estimada e  $\epsilon$  representando o erro de medições com densidade  $P_\epsilon(\epsilon)$ . usando as técnicas de propagação da densidade, temos que a função de verossimilhança é dada por

$$p(y|\theta) = P_\epsilon(y - \theta^2)$$

se  $\epsilon \sim N(0, \sigma_\epsilon^2)$  entao

$$p(y|\theta) = \frac{1}{\sqrt{2\pi \cdot \sigma_\epsilon^2}} \cdot \exp \left\{ -\frac{1}{2} \cdot \frac{(y - \theta^2)^2}{\sigma_\epsilon^2} \right\}$$

se dispomos de  $N$  medidas, elas são iid, e

$$p(y|\theta) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi \cdot \sigma_\epsilon^2}} \cdot \exp \left\{ -\frac{1}{2} \cdot \frac{(y_n - \theta^2)^2}{\sigma_\epsilon^2} \right\}$$

$$\text{e } \ln(p(y|\theta)) = -N \cdot \ln(\sqrt{2\pi \cdot \sigma_\epsilon^2}) - \frac{1}{2 \cdot \sigma_\epsilon^2} \cdot \sum_{n=1}^N (y_n - \theta^2)^2$$

O máximo de  $\ln(p(y|\theta))$  pertence a

$$\begin{aligned} \frac{\partial \ln(p(y|\theta))}{\partial \theta} \Big|_{\theta=\hat{\theta}} &= \frac{1}{2 \cdot \sigma_\epsilon^2} \cdot \sum_{n=1}^N 2 \cdot (y_n - \hat{\theta}^2) \cdot 2 \cdot \hat{\theta} \\ &= \frac{2}{\sigma_\epsilon^2} \cdot \sum_{n=1}^N (y_n - \hat{\theta}^2) \cdot \hat{\theta} \end{aligned}$$

$$b) \hat{\theta}_2 = \sqrt{\bar{y}} \Rightarrow \frac{\partial^2 \ln(p(y|\theta))}{\partial \theta^2} = -\frac{4N}{\sigma_E^2} \cdot \bar{y}$$

$\hat{\theta}_2$  é máxima se  $\bar{y} > 0$

$$c) \hat{\theta}_3 = -\sqrt{\bar{y}} \Rightarrow \frac{\partial^2 \ln(p(y|\theta))}{\partial \theta^2} = -\frac{4N}{\sigma_E^2} \cdot \bar{y}$$

$\hat{\theta}_3$  é máxima se  $\bar{y} > 0$ .

Portanto, sendo  $\bar{y}$  um número real, temos como estimativas:

$$\hat{\theta} = \begin{cases} 0, & \text{se } \bar{y} < 0 \\ \pm\sqrt{\bar{y}}, & \text{se } \bar{y} \geq 0 \end{cases}$$

Observa-se que para  $y \geq 0$ ,  $p(y|\hat{\theta}_1) = p(y|\hat{\theta}_2)$ , ou seja, o problema de estimações MV tem duas soluções que correspondem aos modos de  $p(y|\theta)$ .

Deve ser observado que, embora o parâmetro  $\theta$  a ser estimado tenha sido até então considerado uma constante, a estimativa  $\hat{\theta}$  é uma variável aleatória. Por exemplo, no caso anterior,  $\hat{\theta} = \sqrt{\bar{y}}$  é uma vauma vez que  $\hat{\theta}$  é uma função da va  $\bar{y}$ , que por sua vez é função das vas  $y_1, \dots, y_N$ . Sendo assim  $\hat{\theta}$  possui uma função de densidade associada.

E ainda,  $\hat{\theta}$  possui uma variância (medida de <sup>(108)</sup>  
incerteza). Como  $\hat{\theta} = \sqrt{g}$ , sua variância pode ser  
obtida da variância de  $\bar{y}$

$$\text{Var}\{\bar{y}\} = \text{Var}\left\{ \frac{1}{N} \cdot \sum_{n=1}^N y_n \right\}$$

Sendo  $y_n = \theta^2 + \epsilon_n$ , com  $\epsilon_n \sim N(0, \sigma_\epsilon^2)$ , entao

$\text{Var}\{y_n\} = \sigma_\epsilon^2$ . Sendo  $\bar{y}$  uma função linear de  $y_1, y_2, \dots, y_N$ , entao

$$\begin{aligned} \text{Var}\{\bar{y}\} &= \frac{1}{N^2} \cdot \text{Var}\{y_1\} + \dots + \frac{1}{N^2} \cdot \text{Var}\{y_N\} \\ &= \frac{1}{N^2} \cdot \sum_{n=1}^N \text{Var}\{y_n\} = \frac{1}{N^2} \cdot \sum_{n=1}^N \sigma_\epsilon^2 = \frac{1}{N} \sigma_\epsilon^2 \\ \sigma_{\bar{y}}^2 &= \frac{\sigma_\epsilon^2}{N} \end{aligned}$$

Então,  $\text{Var}\{\hat{\theta}\}$  pode ser obtido através da variância  
de  $\bar{y}$  por meio da propagação de covariâncias.

Considerando agora a teoria em geral, se  $\hat{\theta}$  é  
uma estimativa de  $\theta$  ( $\hat{\theta}^*$  não é obviamente de  
variância mínima), a relação

$$\text{Var}\{\hat{\theta}^*\} \geq \left\{ E \left\{ \left( \frac{\partial \ln(p(y|\theta))}{\partial \theta} \right)^2 \right\} \right\}^{-1}$$

com o turno à direita podendo ser substituída por

$$\left\{ E \left\{ \left( \frac{\partial \ln(p(y|\theta))}{\partial \theta} \right)^2 \right\} \right\}^{-1} = - \left\{ E \left\{ \frac{\partial^2 \ln(p(y|\theta))}{\partial \theta^2} \right\} \right\}^{-1}$$

é conhecida como função logarítmica de Cramer-Rao.  
Este resultado implica se  $\hat{\theta}$  for uma estimativa  
máx fundamental de  $\theta$ . É saido  $\hat{\theta}$  uma estimativa  
máx fundamental, e de variância mínima,

$$\frac{\partial \ln(p(y|\theta))}{\partial \theta} = (\hat{\theta}(y) - \theta) \cdot c(\theta) = 0,$$

com  $c(\theta)$  sendo uma função independente de  $y$   
(dados coletados). O limite inferior de Cramer-Rao  
é alcançado por  $\theta^* = \hat{\theta}$ . Observa-se que a relação  
acima tem como solução  $c(\theta)=0$  e  $\hat{\theta}(y)=\theta$ .  
 $c(\theta)=0$  é independente das medições, sendo assim  
desinteressante para nós. A única solução interessante  
é obtida por  $\hat{\theta}(y)=\theta$ .

Exemplo: Estimação não-linear (pg 105) revisada.

Temos que

$$\begin{aligned}\frac{\partial \ln(p(y|\theta))}{\partial \theta} &= \frac{2}{\sigma_e^2} \sum_{n=1}^N (y_n - \hat{\theta}^2) \cdot \hat{\theta} \\ &= \frac{2}{\sigma_e^2} \left\{ N \bar{y} - N \hat{\theta}^2 \right\} \cdot \hat{\theta} \\ &= \underbrace{\frac{2N}{\sigma_e^2} \cdot \hat{\theta}}_{c(\theta)|_{\theta=\hat{\theta}}} \cdot \underbrace{\left\{ \bar{y} - \hat{\theta}^2 \right\}}_{(\hat{\theta}(y)-\theta)}\end{aligned}$$

Então, a solução  $\hat{\theta}=0$  não é interessante pois  
não depende de  $y_1, \dots, y_N$

## 2.15.3 Estimacões Bayesiana

No paradigma Bayesiano,  $\theta$  é tratado como uma variável distribuída a priori dada por  $P_\theta(\theta)$ . Com a realizações de medidos  $y$  (i.e., vetor de medidos), então, pela regra de Bayes,

$$P(\theta|y) = \frac{P(y|\theta) \cdot P_\theta(\theta)}{P_y(y)} = \frac{P(y|\theta) \cdot P_\theta(\theta)}{\int_{-\infty}^{\infty} P(y|\theta) \cdot P_\theta(\theta) d\theta}$$

Na relação acima temos  $P(y|\theta)$  que é a função de verossimilhança.  $P(\theta|y)$  é chamada de densidade a posteriori. Sendo  $P_y(y)$  independente de  $\theta$ ,  $P(\theta|y)$  é também escrita como

$$P(\theta|y) = \beta \cdot P(y|\theta) \cdot P_\theta(\theta) \quad \left. \begin{array}{l} \\ \end{array} \right\} \beta = \frac{1}{P_y(y)}$$

$\propto P(y|\theta) \cdot P_\theta(\theta)$   
o símbolo "proporcional a"

Ao assumir independência entre os componentes  $y_n$  do vetor  $y$ , usa-se  $P(y|\theta) = \prod_{n=1}^N P(y_n|\theta)$ , o que facilita enormemente a determinação de  $P(\theta|y)$ .

No caso discrete, a regra de Bayes se escreve como

$$F(\theta_i|y_j) = \frac{F(y_j|\theta_i) \cdot F_\theta(\theta_i)}{F_y(y_j)} = \frac{F(y_j|\theta_i) \cdot F_\theta(\theta_i)}{\sum_{k=1}^{N_0} F(y_j|\theta_k) \cdot F_\theta(\theta_k)}$$

Na estimacões bayesiana, dois tipos de estimativas são comumente empregadas: a média  $E\{\theta|y\}$  e o máxímo a posteriori (MAP)  $\arg \max_{\theta} P(\theta|y)$ .

## (11)

### 2.15.3.1 Média como estimativa Bayesiana

Nesta abordagem,

$$\hat{\theta} = E\{\theta|y\} = \int_{-\infty}^{\infty} \theta \cdot p(\theta|y) \cdot d\theta$$

que pode ser determinada de forma analítica ou através de simulações  $\theta_1, \dots, \theta_m$  usando Monte Carlo:

$$\hat{\theta} \approx \frac{1}{m} \cdot \sum_{k=1}^m \theta_k$$

com as amostras  $\theta_k$  geradas a partir de  $p(\theta|y)$ .

Considerando  $\theta$  o verdadeiro parâmetro, temos

$$Var\{\hat{\theta}|y\} = E\{(\hat{\theta} - E\{\hat{\theta}\})^2 | y\}$$

como variância da estimativa  $\hat{\theta}$ . Se  $E\{\hat{\theta}\} = \theta$ , então

$$\begin{aligned} Var\{\hat{\theta}|y\} &= E\{(\hat{\theta} - \theta)^2 | y\} \\ &= \int_{-\infty}^{\infty} (\hat{\theta} - \theta)^2 \cdot p(\theta|y) d\theta \end{aligned}$$

Seu  $Var\{\hat{\theta}|y\} \geq 0$ , então esta é uma função cujo mínimo satisfaaz a

$$\frac{\partial Var\{\hat{\theta}|y\}}{\partial \hat{\theta}} = 2 \cdot \int_{-\infty}^{\infty} (\hat{\theta} - \theta) \cdot p(\theta|y) d\theta = 0$$

Então,

$$\underbrace{\hat{\theta} \cdot \int_{-\infty}^{\infty} p(\theta|y) d\theta}_{1} = \int_{-\infty}^{\infty} \theta \cdot p(\theta|y) d\theta$$

Assim

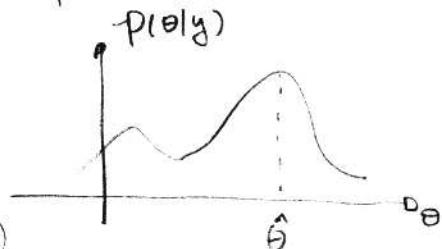
$$\hat{\theta} = \int_{-\infty}^{\infty} \theta \cdot p(\theta|y) d\theta = E\{\theta|y\}$$

Portanto,  $\hat{\theta} = E\{\theta|y\}$  é a estimativa da variância mínima de  $\theta$ . Observe que  $\hat{\theta}$  é uma va pois depende de  $y$ , que é uma va.

### 2.15.3.2 MAP como estimativa Bayesiana

A estimativa MAP é dada por

$$\hat{\theta} = \arg \max_{\theta} p(\theta|y) \Rightarrow$$



Ou seja,  $\hat{\theta}$  é o modo principal de  $p(\theta|y)$  e portanto a

$$\left. \frac{\partial p(\theta|y)}{\partial \theta} \right|_{\theta=\hat{\theta}} = 0 \quad \text{e} \quad \left. \frac{\partial^2 p(\theta|y)}{\partial \theta^2} \right|_{\theta=\hat{\theta}} < 0$$

Sendo  $p(\theta|y) = \frac{p(y|\theta) \cdot p_\theta(\theta)}{p_y(y)} = \beta \cdot p(y|\theta) \cdot p_\theta(\theta)$

então  $\frac{\partial p(\theta|y)}{\partial \theta} = \beta \cdot \frac{\partial [p(y|\theta) \cdot p_\theta(\theta)]}{\partial \theta}$ , com  $\beta = p_y(y)$

O que significa que  $p_y(y)$  é desnecessário para calcular a estimativa MAP pois esta portanto a

$$\frac{\partial p(\theta|y)}{\partial \theta} = 0 \Rightarrow \frac{\partial [p(y|\theta) \cdot p_\theta(\theta)]}{\partial \theta} = 0$$

sendo  $\beta \neq 0$ .

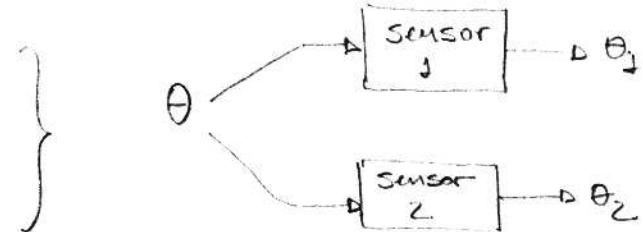
## Exemplo: Fusão de dados multivariacionais

(113)

Considere dois sensores fornecendo medidas  $\theta_1$  e  $\theta_2$  da grandeza  $\theta$  de acordo com os modelos

$$\theta_1 = \theta + \varepsilon_1$$

$$\theta_2 = \theta + \varepsilon_2$$



com  $\varepsilon_1 \sim N(0, \sigma_{\varepsilon_1}^2)$  e  $\varepsilon_2 \sim N(0, \sigma_{\varepsilon_2}^2)$ . Temos que

$$E\{\theta_1\} = E\{\theta_2\} = \theta \quad \text{e} \quad \theta_1 \sim N(\theta, \sigma_{\varepsilon_1}^2) \quad \text{e} \quad \theta_2 \sim N(\theta, \sigma_{\varepsilon_2}^2).$$

Se nenhum conhecimento a priori de  $\theta$  é dado, entao usar  $\hat{\theta} = \theta_1$  como estimativa inicial permite construir esta informação a priori. Assim temos,  $\hat{\theta} \sim N(\theta_1, \sigma_{\varepsilon_1}^2)$

e

$$P_{\theta_1}(\theta) = \frac{1}{(2\pi\sigma_{\varepsilon_1}^2)^{1/2}} \cdot \exp \left\{ -\frac{1}{2} \cdot \frac{(\theta - \theta_1)^2}{\sigma_{\varepsilon_1}^2} \right\}$$

Entao,  $P_\theta(\theta) = P_{\theta_1}(\theta)$  entra como informação a priori.

$$E \quad P(\theta_2 | \theta) = \frac{1}{(2\pi\sigma_{\varepsilon_2}^2)^{1/2}} \cdot \exp \left\{ -\frac{1}{2} \cdot \frac{(\theta_2 - \theta)^2}{\sigma_{\varepsilon_2}^2} \right\}$$

Assim,

$$P(\theta_2 | \theta) \cdot P_\theta(\theta) = \frac{1}{(4\pi^2 \sigma_{\varepsilon_1}^2 \cdot \sigma_{\varepsilon_2}^2)^{1/2}} \cdot \exp \left\{ -\frac{1}{2} \cdot \left\{ \frac{(\theta_2 - \theta)^2}{\sigma_{\varepsilon_2}^2} + \frac{(\theta - \theta_1)^2}{\sigma_{\varepsilon_1}^2} \right\} \right\}$$

Como

$$\hat{\theta} = \arg \max_{\theta} P(\theta_2 | \theta) \cdot P_\theta(\theta) = \arg \max_{\theta} \ln (P(\theta_2 | \theta) \cdot P_\theta(\theta))$$

(114)

$$\ln(P(\theta_2|\theta) \cdot p_\theta(\theta)) = -\ln(2\pi\sigma_{\epsilon_j}^2\sigma_{\epsilon_2}^2) - \frac{1}{2} \cdot \left\{ \frac{(\theta_2-\theta)^2}{\sigma_{\epsilon_2}^2} + \frac{(\theta-\theta_j)^2}{\sigma_{\epsilon_j}^2} \right\}$$

e então,

$$\frac{\partial \ln(P(\theta_2|\theta) \cdot p_\theta(\theta))}{\partial \theta} \Big|_{\theta=\hat{\theta}} = -\frac{1}{2} \cdot \left\{ \frac{2 \cdot (\theta_2 - \hat{\theta})(-1)}{\sigma_{\epsilon_2}^2} + \frac{2 \cdot (\hat{\theta} - \theta_j)}{\sigma_{\epsilon_j}^2} \right\} = 0$$

Portanto,

$$\frac{\theta_2 - \hat{\theta}}{\sigma_{\epsilon_2}^2} = \frac{\hat{\theta} - \theta_j}{\sigma_{\epsilon_j}^2}$$

$$\theta_2 \cdot \sigma_{\epsilon_j}^2 - \hat{\theta} \cdot \sigma_{\epsilon_j}^2 = \hat{\theta} \cdot \sigma_{\epsilon_2}^2 - \theta_j \cdot \sigma_{\epsilon_2}^2$$

$$\hat{\theta} \cdot (\sigma_{\epsilon_j}^2 + \sigma_{\epsilon_2}^2) = \theta_j \cdot \sigma_{\epsilon_2}^2 + \theta_2 \cdot \sigma_{\epsilon_j}^2$$

$$\hat{\theta} = \frac{\sigma_{\epsilon_2}^2}{\sigma_{\epsilon_j}^2 + \sigma_{\epsilon_2}^2} \cdot \theta_j + \frac{\sigma_{\epsilon_j}^2}{\sigma_{\epsilon_j}^2 + \sigma_{\epsilon_2}^2} \cdot \theta_2$$

Portanto,  $\hat{\theta}$  é uma média ponderada de  $\theta_j$  e  $\theta_2$  com pesos  $\alpha = \sigma_{\epsilon_2}^2 / (\sigma_{\epsilon_j}^2 + \sigma_{\epsilon_2}^2)$  e  $1 - \alpha$  talis que

$$\hat{\theta} = \alpha \cdot \theta_j + (1 - \alpha) \cdot \theta_2$$

A variância de  $\hat{\theta}$  é dada por

$$\sigma_{\hat{\theta}}^2 = \alpha^2 \cdot \sigma_{\epsilon_j}^2 + (1 - \alpha)^2 \cdot \sigma_{\epsilon_2}^2$$

Pode ser verificado que, dentro a classe de estimadores  $\hat{\theta} = \alpha \cdot \theta_j + (1 - \alpha) \cdot \theta_2$ , com  $0 \leq \alpha \leq 1$ ,

a variância mínima é alcançada com

$$\alpha_{\text{MV}} = \sigma_{\epsilon_2}^2 / (\sigma_{\epsilon_1}^2 + \sigma_{\epsilon_2}^2)$$

que satisfaz a

$$\frac{\partial \hat{\theta}^2}{\partial \alpha} \Big|_{\alpha=\alpha_{\text{MV}}} = 0.$$

E ainda,  $\hat{\theta}$  pode ser escrito como

$$\hat{\theta} = \theta_j + \underbrace{(1-\alpha_{\text{MV}}) \cdot (\theta_2 - \theta_j)}_{g_{\text{MV}}}$$

com  $g_{\text{MV}}$  sendo o ganho que resulta na variância mínima de  $\hat{\theta}$  obtida de  $\alpha_{\text{MV}}$  e que multiplica a diferença entre  $\theta_2$  e  $\theta_j$ .  $(\theta_2 - \theta_j)$  é chamado de inovações,  $\theta_2$  é a medicaçā e  $\theta_j$  é a estimativa a priori de  $\theta$ . Observar que

$$g_{\text{MV}} = \frac{\sigma_{\epsilon_1}^2}{\sigma_{\epsilon_2}^2 + \sigma_{\epsilon_1}^2} = \begin{cases} 0, & \text{se } \sigma_{\epsilon_1}^2 = 0 \text{ e } \sigma_{\epsilon_2}^2 \neq 0 \\ 1, & \text{se } \sigma_{\epsilon_1}^2 \neq 0 \text{ e } \sigma_{\epsilon_2}^2 = 0 \end{cases}$$

Assim, se a informação a priori é precisa ( $\sigma_{\epsilon_1}^2 = 0$ ), temos  $g_{\text{MV}} = 0$  e  $\hat{\theta} = \theta_j$ . Ou seja, o estimador "contia" apenas na informação a priori. Seus a medicas precisa ( $\sigma_{\epsilon_2}^2 = 0$ ),  $g_{\text{MV}} = 1$  e  $\hat{\theta} = \theta_2$ , prevalecendo como estimativa  $\theta_2$ .

(116)

Neste exemplo específico,  $g_m$  é chamado de Ganso de Kalman, pois o problema é linear, Gaussiano, e o filtro de Kalman, se aplicado, dará o mesmo resultado. Assim, nestas condições, o filtro de Kalman é um caso especial do estimador MAP Bayesiano (para sistemas lineares Gaussianos) → continuo em 116a.

Exemplo: Problema não-linear (pg 105) em versão

Bayesiana:  $y = \theta^2 + \epsilon$ ,  $\epsilon \sim N(0, \sigma_\epsilon^2)$

com  $\theta \in \mathbb{R}$ . Dadas medições  $y_1, \dots, y_N$ , foi obtida

$$\begin{aligned} p(y|\theta) &= \prod_{n=1}^N \frac{1}{(2\pi\sigma_\epsilon^2)^{1/2}} \cdot \exp \left\{ -\frac{1}{2} \cdot \frac{(y_n - \theta^2)^2}{\sigma_\epsilon^2} \right\} \\ &= (2\pi\sigma_\epsilon^2)^{-N/2} \cdot \exp \left\{ \sum_{n=1}^N -\frac{1}{2} \frac{(y_n - \theta^2)^2}{\sigma_\epsilon^2} \right\} \end{aligned}$$

Se tivermos como conhecimento a priori  $\theta \sim N(\theta_0, \sigma_{\theta_0}^2)$  entao  $p_\theta(\theta) = \frac{1}{(2\pi\sigma_{\theta_0}^2)^{1/2}} \cdot \exp \left\{ -\frac{1}{2} \frac{(\theta - \theta_0)^2}{\sigma_{\theta_0}^2} \right\}$

Portanto,

$$p(y|\theta) \cdot p_\theta(\theta) = (2\pi\sigma_\epsilon^2)^{-N/2} \cdot (2\pi\sigma_{\theta_0}^2)^{-1/2} \cdot \exp \left\{ -\frac{1}{2} \left\{ \frac{(\theta - \theta_0)^2}{\sigma_{\theta_0}^2} + \sum_{n=1}^N \frac{(y_n - \theta^2)^2}{\sigma_\epsilon^2} \right\} \right\}$$

$$\begin{aligned} \ln(p(y|\theta) \cdot p_\theta(\theta)) &= -\frac{N}{2} \cdot \ln(2\pi\sigma_\epsilon^2) - \frac{1}{2} \cdot \ln(2\pi\sigma_{\theta_0}^2) \\ &\quad - \frac{1}{2} \cdot \left\{ \frac{(\theta - \theta_0)^2}{\sigma_{\theta_0}^2} + \sum_{n=1}^N \frac{(y_n - \theta^2)^2}{\sigma_\epsilon^2} \right\} \end{aligned}$$

→ continua em 117.

(116a)

A variância da  $\hat{\theta}$  considerando  $\alpha_{mv}$  é dada por

$$\sigma_{\hat{\theta}}^2 = \alpha_{mv}^2 \cdot \sigma_{\epsilon_2}^2 + (1-\alpha_{mv})^2 \cdot \sigma_{\epsilon_1}^2$$

com  $\alpha_{mv} = \frac{\sigma_{\epsilon_2}^2}{\sigma_{\epsilon_1}^2 + \sigma_{\epsilon_2}^2}$ , temos

$$\sigma_{\hat{\theta}}^2 = \frac{\sigma_{\epsilon_2}^4 \cdot \sigma_{\epsilon_2}^2}{(\sigma_{\epsilon_1}^2 + \sigma_{\epsilon_2}^2)^2} + \frac{\sigma_{\epsilon_1}^4 \cdot \sigma_{\epsilon_2}^2}{(\sigma_{\epsilon_1}^2 + \sigma_{\epsilon_2}^2)^2}$$

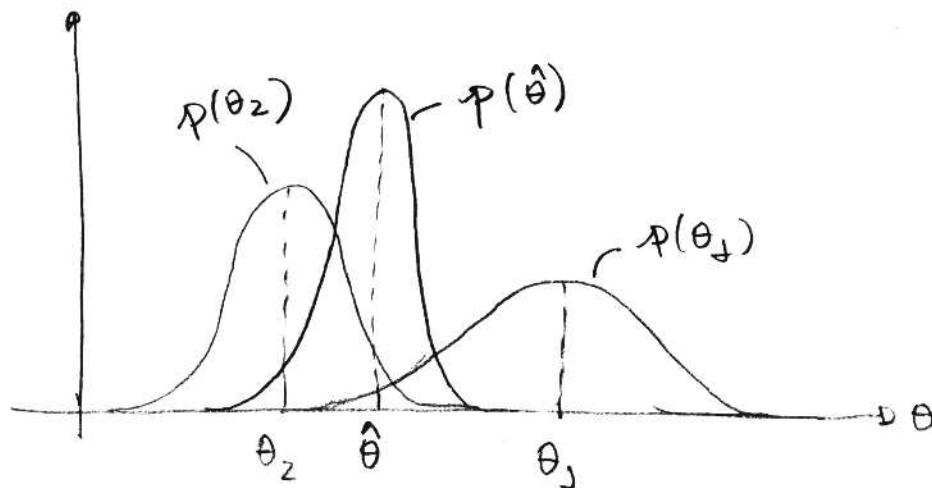
$$\sigma_{\hat{\theta}}^2 = \frac{\sigma_{\epsilon_1}^2 \cdot \sigma_{\epsilon_2}^2 \cdot (\sigma_{\epsilon_2}^2 + \sigma_{\epsilon_1}^2)}{(\sigma_{\epsilon_1}^2 + \sigma_{\epsilon_2}^2)^2} = \frac{\sigma_{\epsilon_1}^2 \cdot \sigma_{\epsilon_2}^2}{(\sigma_{\epsilon_1}^2 + \sigma_{\epsilon_2}^2)}$$

Pode ser verificado que

$$\sigma_{\hat{\theta}}^2 \leq \sigma_{\epsilon_2}^2 \quad \text{e} \quad \sigma_{\hat{\theta}}^2 \leq \sigma_{\epsilon_1}^2$$

ou seja, a incerteza associada a  $\hat{\theta}$  é menor ou igual à menor das incertezas de  $\theta_1$  e  $\theta_2$  (medidas).

Isso justifica o uso de técnicas estatísticas em fusões de dados sensoriais. Percebe-se que



$\in \hat{\theta}$  satisfaz a

$$\frac{2 \cdot (\hat{\theta} - \theta_0)}{\sigma_{\theta_0}^2} + \sum_{n=1}^N \frac{2 \cdot (y_n - \hat{\theta}^2) \cdot (-2) \cdot \hat{\theta}}{\sigma_e^2} = 0$$

$$\frac{(\hat{\theta} - \theta_0)}{\sigma_{\theta_0}^2} - \frac{2}{\sigma_e^2} \cdot \left\{ \underbrace{\sum_{n=1}^N y_n \cdot \hat{\theta}}_{\hat{\theta} \cdot N \bar{y}} - \underbrace{\sum_{n=1}^N \hat{\theta}^3}_{\hat{\theta}^3 \cdot N} \right\} = 0$$

$$\frac{\hat{\theta} - \theta_0}{\sigma_{\theta_0}^2} - \frac{2 \cdot \hat{\theta} \cdot N \bar{y}}{\sigma_e^2} + \frac{2 \cdot \hat{\theta}^3 \cdot N}{\sigma_e^2} = 0$$

$$\frac{\hat{\theta}}{\sigma_{\theta_0}^2} - \frac{2 \cdot \hat{\theta} \cdot N \bar{y}}{\sigma_e^2} + \frac{2 \cdot \hat{\theta}^3 \cdot N}{\sigma_e^2} = \frac{\theta_0}{\sigma_{\theta_0}^2}$$

$$\boxed{\hat{\theta}^3 \cdot \frac{N \cdot 2}{\sigma_e^2} + \hat{\theta} \cdot \left( \frac{1}{\sigma_{\theta_0}^2} - \frac{N \bar{y} \cdot 2}{\sigma_e^2} \right) - \frac{\theta_0}{\sigma_{\theta_0}^2} = 0}$$

Dive ser resolvida esta equação do terceiro grau.

Neste caso, apesar da complexidade, a solução para a estimativa MAP é ainda tratável por seu analítico. Entretanto, por fatores diversos (e.g., múltiplos modos ou casos, multidimensionalidade de  $\theta$ , etc.), a solução analítica é difícil de se obter. Nesses casos, a abordagem computacional é bem-vinda.

(ver [simulações estimativas\\_bayesiana](#))

### 2.15.3.3 Aspectos computacionais

Dadas a equação de Bayes

$$P(\theta|y) = \frac{P(y|\theta) \cdot P(\theta)}{P(y)}$$

foram apresentados dois estimadores que consideram informação a priori: A média condicional  $E\{\theta|y\}$  e o máximo a posteriori (MAP). Apesar da forma simples que  $P(y|\theta)$  e  $P(\theta)$  podem ter, determinar estimativos usando  $p(\theta|y)$  pode ser bastante complicado. E ainda, com o aumento do número de dimensões de  $\theta$ , o uso de técnicas numéricas tradicionais pode ser proibitivo para determinar

$$P(y) = \int_{-\infty}^{\infty} p(y|\theta) \cdot p(\theta) \cdot d\theta$$

$$E\{\theta|y\} = \int_{-\infty}^{\infty} \theta \cdot p(\theta|y) d\theta$$

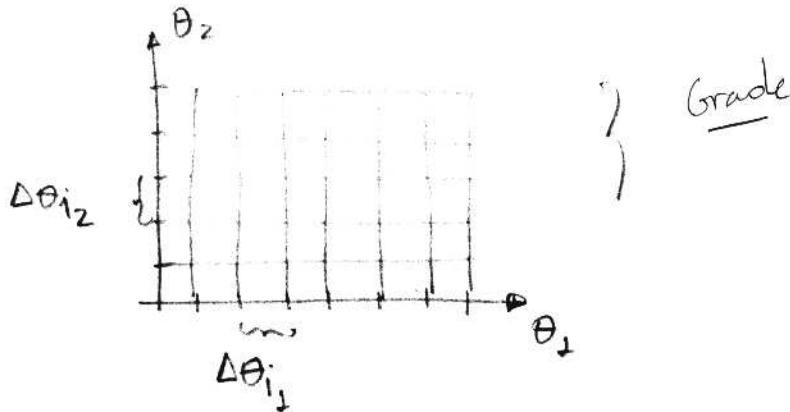
Neste sentido, alguns aproximações numéricas se fazem necessárias.

a) Discretizações do espaço de parâmetros (Grid method) <sup>ou lattice</sup>

Uma abordagem clássica na determinação de  $p(\theta|y)$  consiste em discretizar o espaço do parâmetro  $\theta$  em  $\theta_i = [\theta_{i1}, \dots, \theta_{in}]$ , com

$$\theta_{ij} - \theta_{ij-1} = \Delta\theta_{ij}, \dots, \theta_{in} - \theta_{in-1} = \Delta\theta_{in}.$$

Por exemplo, com  $n=2$ :



Esta discretização somente faz sentido se os parâmetros  $\theta_1, \dots, \theta_n$  tiverem valores limitados a intervalos. Assim,  $p(\theta|y)$  é estimada para cada  $\theta$  na grade. Diversas etapas envolvem:

$$p(y, \theta_i) = p(y|\theta_i) \cdot p(\theta_i)$$

e calcular

$$p(y) = \sum_{i_1} \cdots \sum_{i_n} \underbrace{p(y|\theta_i)}_{p(y,\theta_i)} \cdot p(\theta_i) \cdot \Delta i_1 \cdots \Delta i_n$$

De posso da  $p(y)$  e  $p(y, \theta_i)$ , temos outras

$$p(\theta_i|y) = \frac{p(y, \theta_i)}{p(y)} \quad (\text{IS})$$

b) Amostragem de importância (Importance Sampling)

A ideia principal do IS é aproximar integrais do tipo  $\int f(x)dx$  usando uma função de densidade auxiliar  $g(x)$  de onde amostras são geradas:

$$\begin{aligned} \int f(x)dx &= \int \left( \frac{f(x)}{g(x)} \right) \cdot g(x) \cdot dx \\ &= \frac{1}{M} \cdot \sum_{m=1}^M \frac{f(x_m)}{g(x_m)} \end{aligned}$$

com  $x_m \sim g(x)$ .

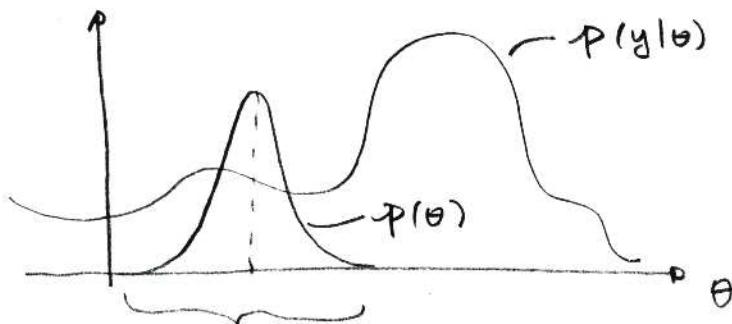
$g(\theta)$  é chamada de densidade de importância.  
De forma similar aos métodos acitação/rejeição,  
 $g(\theta)$  deve ter a forma mais próxima de  $p(\theta)$ , de  
forma que a estimativa da integral seja melhor  
possível.

O emprego deste método em estimações bayesiana  
está na avaliação das integrais

$$p(y) = \int_{-\infty}^{\infty} \underbrace{p(y|\theta)}_{f(y)} \cdot p(\theta) \cdot d\theta$$

$$E[\theta|y] = \int_{-\infty}^{\infty} \underbrace{\theta \cdot p(\theta|y)}_{f(\theta)} \cdot d\theta$$

Observar que para o caso de  $p(y)$  for escolhida  
 $g(\theta) = p(\theta)$ , entâs caímos no caso da simulação  
de monte-carlo para aproximações de  $p(y)$ . No  
entanto, esta escolha para  $g(\theta)$  pode não ser  
eficiente (e qualemte não é) se  $p(\theta)$  não for  
da mesma forma de  $p(y|\theta)$ :



amostras p(uem concentradas  
nesta zona.

### c) Sampling-importance-resampling (SIR)

Em vez de aproximar integrais como no método IS, o método SIR gera amostras diretamente da  $p(\theta|y)$ . A idéia consiste em gerar amostras seguindo  $p(\theta)$  e modificá-las de modo a obter amostras da  $p(\theta|y)$ . Seus

$$p(\theta|y) = \frac{p(y|\theta) \cdot p(\theta)}{\int_{-\infty}^{\infty} p(y|\theta) \cdot p(\theta) d\theta} = \frac{f(\theta)}{\int_{-\infty}^{\infty} f(\theta) d\theta}$$

com  $f(\theta) = p(y|\theta) \cdot p(\theta)$  sendo uma função. Se existir uma constante  $M$  que satisfaz a

$$f(\theta) \leq M \cdot p(\theta)$$

então o método aceitação/rejeição pode ser usado. Observa-se que esta relação pode ser escrita como

$$p(y|\theta) \cdot p(\theta) \leq M \cdot p(\theta)$$

o que leva a  $p(y|\theta) \leq M$ . Esta relação é satisfeita com  $\theta = \hat{\theta}_{\text{MV}}$  (estimativa de máximo de verosimilhança) todos entãos que se escolher  $M = p(y|\hat{\theta}_{\text{MV}})$ . Neste caso, o problema de máximo de verosimilhança deve ser resolvido primeiramente. Como  $\int_{-\infty}^{\infty} f(\theta) d\theta \geq 0$ , entãos

$$\underbrace{\frac{f(\theta)}{\int_{-\infty}^{\infty} f(\theta) d\theta}}_{p(\theta|y)} \leq \frac{M \cdot p(\theta)}{\int_{-\infty}^{\infty} f(\theta) d\theta}$$

(122)

ou seja,  $f(\theta) \leq M \cdot p(\theta)$  é sempre válida para  
qualquer  $\theta$  que satisfaga a  $p(\theta|y) \leq \frac{M \cdot p(\theta)}{\int_{-\infty}^{\infty} f(\theta) d\theta}$ . como

foi dito, o método acitadas/rejeições obteria quer  
uma amostra  $\theta_i$  de  $p(\theta)$  e acita-la com probabili-  
dade

$$p(\theta|y)$$


---

$$\frac{M \cdot p(\theta)}{\int_{-\infty}^{\infty} f(\theta) d\theta} = p(y|\theta_i) \cdot p(\theta_i)$$

que, como  $p(\theta|y) \cdot \int_{-\infty}^{\infty} f(\theta) d\theta = f(y)$ , esta proba-  
bilidade é igual a

$$\frac{f(\theta)}{M \cdot p(\theta_i)} = \frac{p(y|\theta_i) \cdot p(\theta_i)}{p(y|\hat{\theta}_m) \cdot p(\theta_i)} = \frac{p(y|\theta_i)}{p(y|\hat{\theta}_m)}$$

Assim, amostras  $\theta_i \sim p(\theta|y)$  são quadras. Se  $M$  não  
for conhecido, estas amostras de  $p(\theta|y)$  podem ser  
quadras da seguinte forma: a partir de  $N$  amostras  
 $\theta_1, \dots, \theta_N$  de  $p(\theta)$ , calcula-se probabilidades

$$q_n = \frac{w_n}{\sum_{n=1}^N w_n}, \text{ com } w_n = \frac{f(\theta_n)}{p(\theta_n)} = \underbrace{p(y|\theta_n)}_{\text{verosimilhança}}$$

e escolhe-se como amostra  $\theta^*$  tirado das amostras  
 $\theta_1, \dots, \theta_N$  com probabilidade  $q_n$ . Ou seja, cada  
 $\theta_n$  de amostra tem probabilidade  $q_n$  de ser sorteado.

(123) Pode-se verificar que  $\theta^*$  é uma amostra de

$\cdot \Pr\{\theta|y\}$  através da seguinte relação:

$$\Pr\{\theta^* \leq a\} = \sum_{n=1}^N q_n \cdot u(\theta_n - a) \Rightarrow \text{Distribuições desuas de } \theta_1, \dots, \theta_N$$

com  $u(\theta_n - a) = \begin{cases} 1, & \text{se } \theta_n \leq a \\ 0, & \text{caso contrário} \end{cases}$ . Então,

$$\Pr\{\theta^* \leq a\} = \frac{\sum_{n=1}^N w_n \cdot u(\theta_n - a)}{\sum_{i=1}^N w_i}$$

então, como

$$\lim_{N \rightarrow \infty} \sum_{n=1}^N w_n \cdot u(\theta_n - a) = \lim_{N \rightarrow \infty} \sum_{n=1}^N p(y|\theta_n) \cdot u(\theta_n - a)$$
$$= N \cdot \int_{-\infty}^a p(y|\theta) \cdot p(\theta) \cdot d\theta$$

pois  $\theta$  foi gerado de  $p(\theta)$

$$\lim_{N \rightarrow \infty} \sum_{n=1}^N w_n = \lim_{N \rightarrow \infty} \sum_{n=1}^N p(y|\theta_n)$$
$$= N \cdot \int_{-\infty}^{\infty} p(y|\theta) \cdot p(\theta) \cdot d\theta$$

Assim,

$$\lim_{N \rightarrow \infty} \Pr\{\theta^* \leq a\} = \frac{\int_{-\infty}^a p(y|\theta) \cdot p(\theta) \cdot d\theta}{\int_{-\infty}^{\infty} p(y|\theta) \cdot p(\theta) \cdot d\theta} = \Pr\{\theta \leq a | y\}$$

Então, como a densidade de  $\Pr\{\theta \leq a | y\}$  é  $p(\theta|y)$ ,

então, com  $N \rightarrow \infty$ ,  $\theta^*$  é uma amostra de  $p(\theta|y)$ .

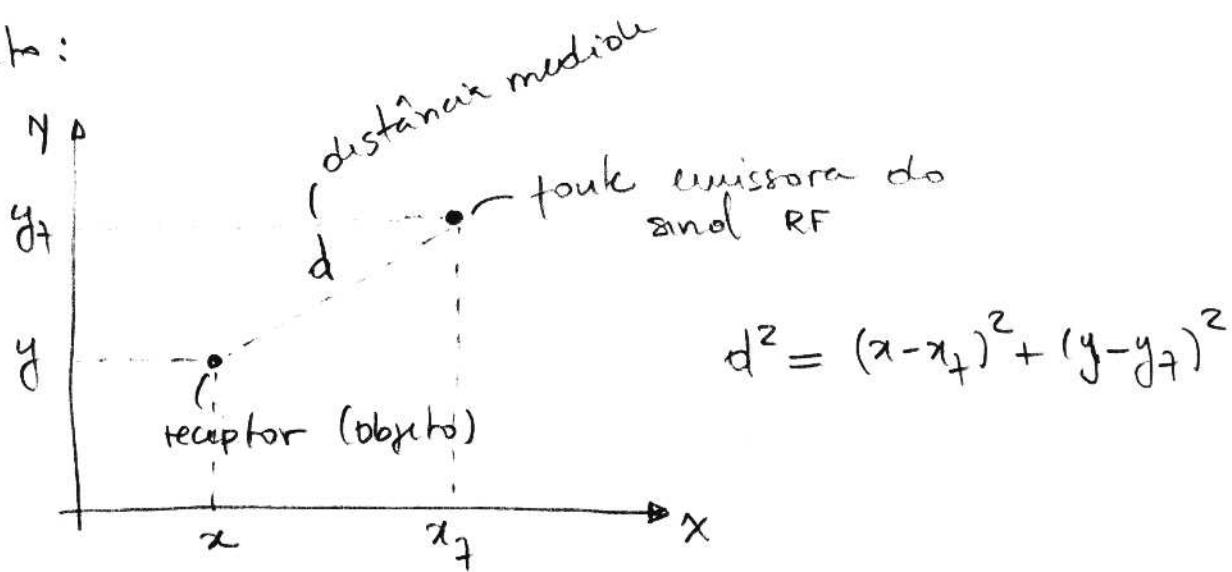
Observa-se que

$$q_n = \frac{p(y|\theta_n)}{\sum_{i=1}^N p(y|\theta_i)}$$

A principal vantagem do algoritmo SIR sobre o método em Grade é que este último se aplica a parâmetros  $\theta$  limitados (e.g., intervalo de interesse) e sua complexidade computacional aumenta com a redução da discretização ( $\Delta\theta_1, \dots, \Delta\theta_N$ ). No SIR, o número de amostras a serem geras de  $p(\theta|y)$  é constante, sendo que este número deve ser suficientemente grande de forma que as amostras possam ser usadas de forma satisfatória para calcular  $\hat{\theta}$ . Além disso,  $p(\theta)$  deve descrever uma distribuição no mínimo pessimista da informação a priori. Caso contrário  $p(\theta)$  não faria parte do suporte de  $p(\theta|y)$ , e poucas amostras utilizárias seriam geradas.

Exemplo: localizações usando RF (posições apena) (125)

Considera o seguinte caso: uma fonte emissora de sinal de RF é usada como "landmark" para a estimativa de posições 2D (no plano) de um objeto (robô, automóvel, pessoa) usando apenas uma medida d da distância do emissor com relação ao objeto:



$(x_f, y_f)$ : coordenadas cartesianas da fonte

$(x, y)$ : " " " do objeto

Definir:

$p(\theta)$ : densidade a priori da posição  $\theta = (x, y)^T$  do objeto, ou ainda  $p(x, y)$ .

$p(d|\theta)$ : densidade condicional de d dado  $\theta$  ou ainda  $p(d|x, y)$

Se modularmos a distância medida d como

$$d = ((x - x_f)^2 + (y - y_f)^2)^{1/2} + \epsilon_d$$

com  $\epsilon_d \sim N(0, \sigma_{\epsilon_d}^2)$

Temos que

$$p(d|\theta) = p_{\epsilon_d} \left( d - \left[ (x-x_f)^2 + (y-y_f)^2 \right]^{1/2} \right)$$

$$= \frac{1}{(2\pi\sigma_{\epsilon_d}^2)^{1/2}} \cdot \exp \left\{ \frac{-1}{2\sigma_{\epsilon_d}^2} \cdot \left( d - \sqrt{(x-x_f)^2 + (y-y_f)^2} \right)^2 \right\}$$

Assim, dados  $p(x,y)$  como informação a priori, entao

$$p(x,y|d) \propto p(d|x,y) \cdot p(x,y)$$

(ver simulações estimativas bayesianas computacionais)