

INVERSÃO LINEAR:

ESTIMADOR DE MÁXIMA VEROSSIMILHANÇA (MAXIMUM LIKELIHOOD)

Os estimadores de Máxima Verossimilhança (EMV) usam informação a priori sobre o comportamento estatístico da variável aleatória (v.a.) que presumivelmente contamina as observações geofísicas (dados). Tais estimadores não usam nem explicita nem implicitamente informação a priori sobre os parâmetros como os estimadores que anteriormente estudamos

Conceitos Fundamentais:

Considere uma v.a. \mathcal{E} sobre a qual conhecemos a função de densidade de probabilidade (f.d.p.) $f(\mathcal{E}, \mu)$ a menos da média μ . Vamos supor que conhecemos N realizações desta v.a., i.e., $\mathcal{E}_1, \dots, \mathcal{E}_N$. Dado uma v.a. \mathcal{E}_i temos associado a ela uma f.d.p. $f(\mathcal{E}_i, \mu)$.

O método Máxima Verossimilhança (MV) consiste em estimar a média μ . Então, no método MV estima-se o valor $\hat{\mu}$ tal que os valores da f.d.p. $f(\mathcal{E}_i, \hat{\mu})$, $i=1, \dots, N$ sejam todos grandes. Em outras palavras o método MV maximiza a sequência particular das realizações $\mathcal{E}_1, \dots, \mathcal{E}_N$.

A função de densidade de probabilidade conjunta das v.a. $\mathcal{E}_1, \dots, \mathcal{E}_N$ é dada por

$$L(\mathbf{\varepsilon}, \mu) = f(\mathcal{E}_1, \mu) \times f(\mathcal{E}_2, \mu) \times \dots \times f(\mathcal{E}_N, \mu),$$

ou seja, a função de densidade de probabilidade conjunta é o produto das N f.d.p. $f(\mathcal{E}_i, \mu)$, $i=1, \dots, N$

$$L(\mathbf{\varepsilon}, \mu) = \prod_{i=1}^N f(\mathcal{E}_i, \mu) \quad (1)$$

O estimador $\hat{\mu}$ da MV da média de $\varepsilon_1, \dots, \varepsilon_N$ será aquele que satisfizer a condição

$$\max_{\mu} L(\boldsymbol{\varepsilon}, \boldsymbol{\mu}) = \max_{\mu} \prod_{i=1}^N f(\varepsilon_i, \mu) \quad (2)$$

Na geofísica as nossas v.a. são os dados observados ($\bar{\mathbf{y}}^o$). Comumente fazemos em cada ponto de medida apenas uma observação, ou seja, em cada ponto temos apenas uma realização de uma v.a.

O estimador $\hat{\mu}$ da MV dos dados observados ($\bar{\mathbf{y}}^o$) será aquele que satisfizer a condição

$$\max_{\mu} L(\bar{\mathbf{y}}^o, \boldsymbol{\mu}) = \max_{\mu} f(\bar{\mathbf{y}}^o, \mu) \quad (3)$$

Mas quem é a média μ dos dados observados $\bar{\mathbf{y}}^o$?

Cálculo da $E\{\bar{\mathbf{y}}^o\}$:

Considerando as premissas estatísticas estabelecidas na Figura 1 e as propriedades da esperança, temos:

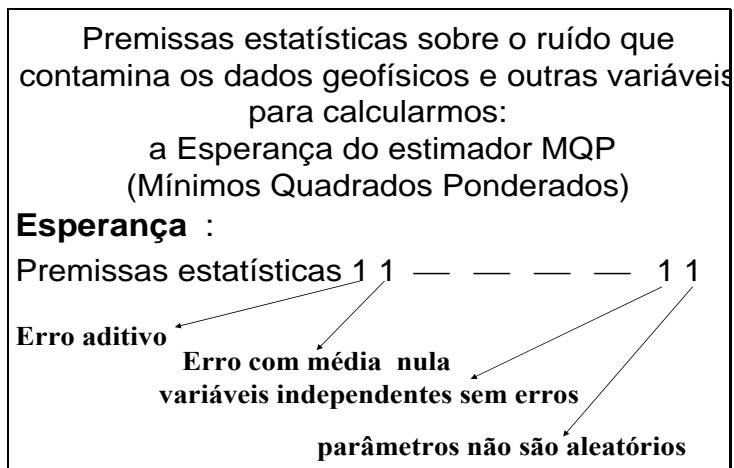


Figura 1

Usando a premissa 1 que os erros são aditivos ($\bar{y}^o = \bar{y}^c + \bar{\varepsilon}$) e usando a informação que em um problema linear a componente determinística (vetor dos

dados ajustados ou calculados) é $\bar{y}^c = \bar{\bar{A}} \bar{p}$ temos que

$$\bar{y}^o = \bar{\bar{A}} \bar{p} + \bar{\varepsilon}$$

Calculando a esperança da v.a. que são os dados (\bar{y}^o) temos

$$E\{\bar{y}^o\} = E\{\bar{\bar{A}} \bar{p}\} + E\{\bar{\varepsilon}\}$$

Usando a premissa 7 que as variáveis independentes não são v.a. e a premissa 8 que os parâmetros não são aleatórios temos

$$E\{\bar{y}^o\} = \bar{\bar{A}} \bar{p} + E\{\bar{\varepsilon}\}$$

Pela premissa 2 os erros tem média nula o que implica $E[\bar{\varepsilon}] = \bar{0}$ então

$$E\{\bar{y}^o\} = \bar{\bar{A}} \bar{p}. \quad (4)$$

O ESTIMADOR DE MÁXIMA VEROSSIMILHANÇA

Substituindo a equação (4) na equação (3) temos o estimador da MV dos dados observados (\bar{y}^o)

$$\max_{\mu} L(\bar{y}^o, E\{\bar{y}^o\}) = \max f(\bar{y}^o, E\{\bar{y}^o\})$$

$$\max_{\mu} L(\bar{y}^o, \bar{\bar{A}} \bar{p}) = \max f(\bar{y}^o, \bar{\bar{A}} \bar{p}) \quad (5)$$

Caso particular: Distribuição Gaussiana

Considerando apenas uma observação (y_i^o) e considerando o caso particular da distribuição Gaussiana temos que a f.d.p para y_i^o é dada por

$$f(y_i^o, y_i^c) = \frac{1}{(2\pi)^{1/2} \sigma_i} e^{-\frac{(y_i^o - y_i^c)^2}{2\sigma_i^2}}$$

em que σ_i é o desvio padrão da i-ésima observação y_i^o

A função de máxima verossimilhança é dada pelo produto das N f.d.p.

$$f(y_i^o, y_i^c), i=1, \dots, N$$

$$L(\bar{y}^o, \bar{\bar{A}} \bar{p}) = \prod_{i=1}^N \frac{1}{(2\pi)^{1/2} \sigma_i} e^{-\frac{(y_i^o - y_i^c)^2}{2\sigma_i^2}}$$

$$L(\bar{y}^o, \bar{\bar{A}} \bar{p}) = \frac{1}{(2\pi)^{N/2} \prod_{i=1}^N \sigma_i} \prod_{i=1}^N e^{-\frac{(y_i^o - y_i^c)^2}{2\sigma_i^2}} \quad (6)$$

A função acima é a função de máxima verossimilhança que notação matricial pode ser reescrita como

$$L(\bar{y}^o, \bar{\bar{A}} \bar{p}) = \frac{1}{(2\pi)^{N/2} \bar{\bar{W}}^{1/2}} e^{-\frac{1}{2} \left\{ (\bar{y}^o - \bar{\bar{A}} \bar{p})^T \bar{\bar{W}}^{-1} (\bar{y}^o - \bar{\bar{A}} \bar{p}) \right\}} \quad (7)$$

em que $\bar{\bar{W}}$ (N x N) é uma matriz diagonal cujo i-ésimo elemento da diagonal é a variância σ_i^2 do i-ésimo elemento do vetor dos resíduos ($\bar{\epsilon}$) então temos

$$\bar{\bar{W}} = \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ & & & \sigma_N^2 \end{bmatrix}_{N \times N}$$

O estimador de máxima verossimilhança dos dados observados (equação 5) consiste em maximizar a função de máxima verossimilhança (equação 7) para estimarmos o vetor de parâmetros $\bar{\mathbf{p}}$, i.e.:

$$\max_{\bar{\mathbf{p}}} L(\bar{\mathbf{y}}^o, \bar{\mathbf{A}}\bar{\mathbf{p}}) = \max_{\bar{\mathbf{p}}} \frac{1}{(2\pi)^{N/2} \bar{\mathbf{W}}^{1/2}} e^{-\frac{1}{2} \left\{ (\bar{\mathbf{y}}^o - \bar{\mathbf{A}}\bar{\mathbf{p}})^T \bar{\mathbf{W}}^{-1} (\bar{\mathbf{y}}^o - \bar{\mathbf{A}}\bar{\mathbf{p}}) \right\}} \quad (8)$$

Para maximizarmos a função acima basta minimizar o expoente da exponencial.

Portanto a condição necessária para que $L(\bar{\mathbf{y}}^o, \bar{\mathbf{A}}\bar{\mathbf{p}})$ tenha um máximo é

$$\begin{aligned} \bar{\nabla}_{\bar{\mathbf{p}}} \left\{ (\bar{\mathbf{y}}^o - \bar{\mathbf{A}}\bar{\mathbf{p}})^T \bar{\mathbf{W}}^{-1} (\bar{\mathbf{y}}^o - \bar{\mathbf{A}}\bar{\mathbf{p}}) \right\} &= \bar{\mathbf{0}} \\ -2 \bar{\mathbf{A}}^T \bar{\mathbf{W}}^{-1} (\bar{\mathbf{y}}^o - \bar{\mathbf{A}}\hat{\bar{\mathbf{p}}}_L) &= \bar{\mathbf{0}} \end{aligned} \quad (9)$$

que resulta no Estimador de Máxima Verossimilhança (EMV):

$$\hat{\bar{\mathbf{p}}}_L = \left(\bar{\mathbf{A}}^T \bar{\mathbf{W}}^{-1} \bar{\mathbf{A}} \right)^{-1} \bar{\mathbf{A}}^T \bar{\mathbf{W}}^{-1} \bar{\mathbf{y}}^o \quad (\text{EMV}) \quad (10)$$

Note que o estimador de Máxima Verossimilhança é a soma dos resíduos ponderados pelo o inverso da variâncias dos erros. Portanto o EMV é igual ao estimador MQP (mínimos quadrados ponderados que apresentamos no tópico 8 deste curso).

$$\min_{\bar{\mathbf{p}} \in F} \left\{ \bar{\boldsymbol{\varepsilon}}^T \bar{\mathbf{W}}^{-1} \bar{\boldsymbol{\varepsilon}} \right\} = \sum_{i=1}^N \frac{\varepsilon_i^2}{\sigma_i^2}$$

Assim, valem as mesmas observações que fizemos para o MQP, como por exemplo, observações mais confiáveis (variância pequena) produzirão resíduo pequeno.

Equivalência entre estimadores que usam o critério de minimização de uma norma e o EMV

Veremos a Equivalência entre estimadores que usam como critério a minimização de uma norma dos resíduos e certos estimadores de Máxima Verossimilhança. Por simplicidade limitaremos ao caso de um único parâmetro.

1) EMV – Distribuição Gaussiana

Considerando apenas uma observação (y_i^o) e considerando o caso particular da distribuição Gaussiana temos que a f.d.p para y_i^o é dada por

$$f(y_i^o, \mu) = \frac{1}{(2\pi)^{1/2} \sigma} e^{-\frac{(y_i^o - \mu)^2}{2\sigma^2}}$$

em que σ é o desvio padrão da observação y_i^o

Considerando um único parâmetro a ser estimado (μ) e um único desvio padrão para todas as observações (σ) a função de máxima verossimilhança é dada pelo produto das N f.d.p. $f(y_i^o, \mu)$, $i=1, \dots, N$

$$L(\bar{y}^o, \bar{\mathbf{A}} \bar{\mathbf{p}}) = \prod_{i=1}^N \frac{1}{(2\pi)^{1/2} \sigma} e^{-\frac{(y_i^o - \mu)^2}{2\sigma^2}}$$

$$L(\bar{y}^o, \bar{\mathbf{A}} \bar{\mathbf{p}}) = \frac{1}{(2\pi)^{N/2} \sigma^N} e^{-\frac{1}{2} \sum_{i=1}^N \left(\frac{y_i^o - \mu}{\sigma} \right)^2}$$

Para maximizarmos a função acima basta minimizar o expoente da exponencial. Portanto o máximo ocorre quando

$$\text{minimiza-se: } \sum_{i=1}^N \left(\frac{y_i^O - \mu}{\sigma} \right)^2.$$

Para minimizarmos esta expressão em relação ao parâmetro desconhecido μ fazemos

$$\frac{\partial}{\partial \mu} \sum_{i=1}^N \left(\frac{y_i^O - \mu}{\sigma} \right)^2 = 0$$

$$\frac{2}{\sigma^2} \sum_{i=1}^N (y_i^O - \hat{\mu})(-1) = 0.$$

e estimamos

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N y_i^O$$

Assim concluímos que o EMV sob a hipótese que os erros tem distribuição Gaussiana:

1.1) é equivalente a minimização da norma L2 dos resíduos, i.e.,

$$\|\bar{\epsilon}\|_2 = \left(\sum_{i=1}^N |\epsilon_i|^2 \right)^{1/2} \quad \text{com } \epsilon_i = y_i^O - \mu. \text{ Note que a}$$

função com raiz quadrada é monotônica de modo que se $\hat{\mu}$ minimiza $\|\bar{\epsilon}\|_2$

minimiza também $\|\bar{\epsilon}\|_2^{1/2}$ e $\|\bar{\epsilon}\|_2^2$.

1.2) baseia-se na **média amostral**

2) EMV – Distribuição de Laplace

Considerando apenas uma observação (y_i^o) e considerando o caso particular da distribuição de Laplace temos que a f.d.p para y_i^o é dada por

$$f(y_i^o, \mu) = \frac{1}{2\sigma} e^{-\frac{1}{\sigma} |y_i^o - \mu|}$$

em que σ é o desvio padrão da observação y_i^o

Considerando um único parâmetro a ser estimado (μ) e um único desvio padrão para todas as observações (σ) a função de máxima verossimilhança é dada pelo produtório das N f.d.p. $f(y_i^o, \mu)$, $i=1, \dots, N$

$$L(\bar{y}^o, \bar{\bar{A}} \bar{p}) = \prod_{i=1}^N \frac{1}{2\sigma} e^{-\frac{1}{\sigma} |y_i^o - \mu|}$$

$$L(\bar{y}^o, \bar{\bar{A}} \bar{p}) = \frac{1}{2\sigma} e^{-\frac{1}{\sigma} \sum_{i=1}^N |y_i^o - \mu|}$$

Para maximizarmos a função acima basta minimizar o expoente da exponencial. Portanto o máximo ocorre quando

$$\text{minimiza-se: } \sum_{i=1}^N |y_i^o - \mu|.$$

Para minimizarmos esta expressão em relação ao parâmetro desconhecido μ fazemos

$$\frac{\partial}{\partial \mu} \sum_{i=1}^N |y_i^O - \mu| = 0$$

$$\sum_{i=1}^N \text{sgn}(y_i^O - \hat{\mu}) = 0.$$

em que $\text{sgn}(y_i^O - \hat{\mu})$ é o sinal do i-ésimo resíduo. Assim concluímos que o EMV

sob a hipótese que os erros têm distribuição de Laplace :

2.1) é equivalente a minimização da norma L1 dos resíduos, i.e.,

$$\|\overline{\boldsymbol{\varepsilon}}\|_1 = \left(\sum_{i=1}^N |\varepsilon_i| \right) \text{ com } \varepsilon_i = y_i^O - \mu.$$

2.2) baseia-se na **mediana amostral**