

## **INVERSÃO LINEAR:**

### **ESTIMADORES VIA NORMA L1**

Neste tópico iremos ver a definição da norma L1 e seu significado físico. Em seguida iremos falar sobre dois estimadores. O primeiro é o estimador de mínimo dos resíduos via norma L1 e segundo é o estimador da variação total (total variation –TV).

#### **Definição da norma $L_n$**

Por definição, a norma  $\ell_n$  de um vetor genérico M-dimensional  $\bar{v}$  é dada por:

$$\|\bar{v}\|_n = \left[ \sum_{i=1}^M |v_i|^n \right]^{1/n} \quad (1)$$

Portanto a norma  $\ell_1$  de um vetor genérico M-dimensional  $\bar{v}$  é dada por:

$$\|\bar{v}\|_1 = \sum_{i=1}^M |v_i| \quad (2)$$

#### **Estimador via a minimização do vetor dos resíduos segundo a norma 1**

Seja um vetor N-dimensional  $\boldsymbol{\varepsilon}$  que representa a diferença entre os dados geofísicos observados  $\bar{y}^o$  e estimados (ou preditos)  $\bar{y}^c$  tal que

$$\boldsymbol{\varepsilon} = \bar{y}^o - \bar{y}^c$$

No tópico 6 deste curso, deduzimos o estimador de mínimos quadrados (MQ) que consiste em minimizar a soma dos quadrados dos resíduos, i.e.,

$$\min_{\bar{\mathbf{p}} \in F} \{Q\} \equiv \min_{\bar{\mathbf{p}} \in F} \{\|\boldsymbol{\varepsilon}\|_2^2\} \equiv \min_{\bar{\mathbf{p}} \in F} \left\| \bar{\mathbf{y}}^o - \bar{\mathbf{A}} \bar{\mathbf{p}} \right\|_2^2 \equiv \min_{\bar{\mathbf{p}} \in F} \left( \bar{\mathbf{y}}^o - \bar{\mathbf{A}} \bar{\mathbf{p}} \right)^T \left( \bar{\mathbf{y}}^o - \bar{\mathbf{A}} \bar{\mathbf{p}} \right)$$

No problema acima minimizamos o vetor dos resíduos  $\boldsymbol{\varepsilon}$  usando o quadrado da norma L2. Isto equivale a minimizarmos a norma L2. Vimos no tópico 6 que a minimização desta função  $Q$  resulta no estimador dos mínimos quadrados (MQ sobredeterminado). Aqui chamaremos este estimador de

$$\hat{\bar{\mathbf{p}}}_{L2} = \left( \bar{\mathbf{A}}^T \bar{\mathbf{A}} \right)^{-1} \bar{\mathbf{A}}^T \bar{\mathbf{y}}^o \quad (\text{estimador MQ via norma L2}) \quad (3)$$

Neste tópico definimos a função  $Q_1$  como a norma  $\ell_1$  dos resíduos:

$$Q_1 = \|\boldsymbol{\varepsilon}\|_1 = \sum_{i=1}^N |\varepsilon_i|. \quad (4)$$

Então iremos minimizar a função  $Q_1$ , ou seja, minimizaremos o vetor dos resíduos  $\boldsymbol{\varepsilon}$  (N x1) usando a norma  $\ell_1$ :

$$\min_{\bar{\mathbf{p}}} \{Q_1\} = \min_{\bar{\mathbf{p}}} \|\boldsymbol{\varepsilon}\|_1 = \min_{\bar{\mathbf{p}}} \sum_{i=1}^N |\varepsilon_i|. \quad (5)$$

Como  $\boldsymbol{\varepsilon} = \bar{\mathbf{y}}^o - \bar{\mathbf{A}} \bar{\mathbf{p}}$  em um problema linear, então a minimização do vetor dos resíduos  $\boldsymbol{\varepsilon}$  (N x1) usando a norma  $\ell_1$ , pode ser escrita como

$$\min_{\bar{\mathbf{p}}} \{Q_1\} = \min_{\bar{\mathbf{p}}} \|\boldsymbol{\varepsilon}\|_1 = \min_{\bar{\mathbf{p}}} \sum_{i=1}^N |\varepsilon_i| = \min_{\bar{\mathbf{p}}} \left| y^o_i - y^c_i \right| \quad (6)$$

em que  $y^c_i$  é a i-ésimo elemento do vetor de dados ajustados (ou calculados ou preditos), i.e,  $\bar{\mathbf{y}}^c = \bar{\mathbf{A}} \bar{\mathbf{p}}$ .

A condição de mínimo é que a derivada da função  $Q_1$  (equação 4) em relação ao vetor de parâmetros  $\bar{\mathbf{p}}$  seja zero. Note que esta função não é diferenciável se um dos elementos de  $\varepsilon_1, \dots, \varepsilon_N$  é zero. Vamos ignorar este fato e computar a derivada em relação ao k-ésimo parâmetro  $p_k$ , nos pontos em que os elementos de  $\boldsymbol{\varepsilon}$  não são zero:

$$\frac{\partial Q_1}{\partial p_k} = \sum_{i=1}^N \frac{\partial |\varepsilon_i|}{\partial p_k} = \sum_{i=1}^N a_{ik} \operatorname{sgn}(\varepsilon_i) \quad (7)$$

em que  $a_{ik}$  é o  $ik$ -ésimo elemento da matriz de sensibilidade  $\bar{\bar{\mathbf{A}}}$  e  $\operatorname{sgn}(\varepsilon_i)$  é o sinal do  $i$ -ésimo resíduo. Note que a derivada da função  $Q_1$  (equação 4) em relação ao vetor de parâmetros  $\bar{\mathbf{p}}$  depende do resíduo  $\boldsymbol{\varepsilon}$ , como  $\boldsymbol{\varepsilon} = \bar{\mathbf{y}}^o - \bar{\mathbf{y}}^c = \bar{\mathbf{y}}^o - \bar{\bar{\mathbf{A}}}\bar{\mathbf{p}}$  então estamos diante de um problema rigorosamente não-linear. Isto porque a derivada da função depende dos parâmetros desconhecidos  $\bar{\mathbf{p}}$  que queremos estimar.

Apesar do problema de minimização dos resíduos segundo a norma 1 ser um rigorosamente não-linear, uma das alternativas é resolver através de um método iterativo dos mínimos quadrados reponderados (IRLS – iteratively reweighted least squares).

Para tanto vamos reescrever o sinal do  $i$ -ésimo resíduo como

$$\operatorname{sgn}(\varepsilon_i) = \frac{\varepsilon_i}{|\varepsilon_i|},$$

Então podemos reescrever a equação 7 como

$$\frac{\partial Q_1}{\partial p_k} = \sum_{i=1}^N \frac{\partial |\varepsilon_i|}{\partial p_k} = \sum_{i=1}^N a_{ik} \operatorname{sgn}(\varepsilon_i) = \sum_{i=1}^N a_{ik} \frac{1}{|\varepsilon_i|} \varepsilon_i \quad (8)$$

Em notação matricial o gradiente da função  $Q_1$  (equação 4) em relação ao vetor de parâmetros  $\bar{\mathbf{p}}$  é:

$$\bar{\nabla}_{\bar{\mathbf{p}}} \{Q_1\} = \bar{\bar{\mathbf{A}}}^T \bar{\bar{\mathbf{W}}} \boldsymbol{\varepsilon} = \bar{\bar{\mathbf{A}}}^T \bar{\bar{\mathbf{W}}} (\bar{\mathbf{y}}^o - \bar{\bar{\mathbf{A}}}\bar{\mathbf{p}}), \quad (9)$$

em que  $\bar{\bar{\mathbf{W}}}$  ( $N \times N$ ) é uma matriz diagonal de pesos cujo  $i$ -ésimo elemento da diagonal é o inverso do valor absoluto do resíduo estimado na iteração anterior.

$$\bar{\bar{\mathbf{W}}} = \begin{bmatrix} 1/|\varepsilon_1| & & & \\ & 1/|\varepsilon_2| & & \\ & & \ddots & \\ & & & 1/|\varepsilon_N| \end{bmatrix}_{N \times N} \quad (9a)$$

A condição de mínimo é que a derivada da função  $Q_1$  (equação 4) em relação ao vetor de parâmetros  $\bar{\mathbf{p}}$  seja zero. Então fazendo

$$\bar{\nabla}_{\bar{\mathbf{p}}} \{Q_1\} = \bar{\mathbf{A}}^T \bar{\mathbf{W}} (\bar{\mathbf{y}}^o - \bar{\mathbf{A}} \bar{\mathbf{p}}) = \bar{\mathbf{0}}$$

temos

$$\bar{\mathbf{A}}^T \bar{\mathbf{W}} \bar{\mathbf{y}}^o - \bar{\mathbf{A}}^T \bar{\mathbf{W}} \bar{\mathbf{A}} \hat{\mathbf{p}}_{L1} = \bar{\mathbf{0}}.$$

Logo chegamos ao sistema de equações:

$$\bar{\mathbf{A}}^T \bar{\mathbf{W}} \bar{\mathbf{A}} \hat{\mathbf{p}}_{L1} = \bar{\mathbf{A}}^T \bar{\mathbf{W}} \bar{\mathbf{y}}^o \quad (10)$$

em que  $\hat{\mathbf{p}}_{L1}$  é o estimador dos parâmetros via a minimização do vetor dos resíduos através da norma 1.

Como a matriz de peso  $\bar{\mathbf{W}}$  no sistema de equação 10 depende dos resíduos (i.e., depende dos parâmetros), então estamos diante de um sistema de equações não lineares. Uma das alternativas é usar o método iterativo dos mínimos quadrados ponderados (IRLS). Então, na k-ésima iteração estimamos um vetor de parâmetros  $\hat{\mathbf{p}}_{L1}^{(k)}$ . O algoritmo começa com a solução de mínimos quadrados sobre determinado (estimador chamado neste tópico de  $\hat{\mathbf{p}}_{L2}$ , veja a equação 3). Então, na iteração  $k = 1$  estimamos:

$$\hat{\mathbf{p}}_{L1}^{(1)} = \hat{\mathbf{p}}_{L2},$$

em seguida, computamos o vetor de resíduos  $\boldsymbol{\varepsilon}^{(1)} = (\bar{\mathbf{y}}^o - \bar{\mathbf{A}} \hat{\mathbf{p}}_{L1}^{(1)})$  e a matriz de pesos  $\bar{\mathbf{W}}^{(1)}$  (equação 9a) na iteração  $k=1$ . Então, ainda na iteração  $k=1$ , resolvemos o sistema de equações (equação 10):

$$\bar{\mathbf{A}}^T \bar{\mathbf{W}}^{(k)} \bar{\mathbf{A}} \hat{\mathbf{p}}_{L1}^{(k+1)} = \bar{\mathbf{A}}^T \bar{\mathbf{W}}^{(k)} \bar{\mathbf{y}}^o$$

para estimarmos  $\hat{\mathbf{p}}_{L1}^{(2)}$  e o vetor de resíduos  $\boldsymbol{\varepsilon}^{(2)} = (\bar{\mathbf{y}}^o - \bar{\mathbf{A}} \hat{\mathbf{p}}_{L1}^{(2)})$  e a matriz de pesos  $\bar{\mathbf{W}}^{(2)}$  (equação 9a) na iteração  $k=1+1$ . O processo é repetido até que a desigualdade

$$\frac{\|\hat{\mathbf{p}}_{L1}^{(k+1)} - \hat{\mathbf{p}}_{L1}^{(k)}\|_2}{1 + \|\hat{\mathbf{p}}_{L1}^{(k+1)}\|_2} \leq \tau$$

em que  $\tau$  é um número pequeno chamado de tolerância.

Este procedimento falha quando se algum resíduo for igual a zero por causa da matriz de pesos (equação 9a). A solução é definir um valor pequeno  $r$  de modo que se em alguma iteração  $k$  o  $i$ -ésimo resíduo  $\varepsilon_i^{(k)}$  for menor que  $r$  então faz-se  $|\varepsilon_i^{(k)}| = r$ , e conseqüentemente teremos que o  $i$ -ésimo elemento da diagonal da matriz de pesos  $\overline{\mathbf{W}}^{(k)}$  é igual a  $w_{ii}^{(k)} = 1/r$ .

### **Qual o significado físico da minimização da Norma 1 do vetor de resíduos**

Ao minimizarmos o vetor dos resíduos segundo a norma 1 em relação ao vetor de parâmetros  $\bar{\mathbf{p}}$  estamos minimizando a função  $Q_1$  (equação 4). Então minimizamos o sinal dos resíduos. Minimizar o sinal dos resíduos significa que podemos ter valores elevados dos resíduos, porém o somatório dos sinais destes resíduos deve ser mínimo. Isto faz com que a minimização da norma L1 dos resíduos despreze dados observados que sejam espúrios (outliers). Dizemos que a minimização da norma L1 dos resíduos é mais robusta no sentido de permitir grandes resíduos (diferença entre os dados observados e os dados ajustados ou preditos). A Figura 1 mostra em pontos pretos um conjunto de observações geofísicas. Note que há um único dado espúrio (outlier). As retas tracejadas coloridas representam retas de ajustes (dados ajustados ou preditos) minimizando-se as normas 2 (reta L2) e 1 (reta L1) do vetor de resíduos  $\boldsymbol{\varepsilon}$ . O resíduo na  $i$ -ésima coordenada é a diferença entre o dado observado (pontos pretos) e o ajustado (retas ajustadas via normas L1 e L2). Note que a reta ajustada via norma L2 é atraída pelo ponto espúrio (outlier). Ao contrário, a minimização da norma L1 dos resíduos permite grandes resíduos

pois minimiza-se o sinal dos resíduos. Logo a reta ajustada via norma L1 NÃO é atraída pelo ponto espúrio (outlier).

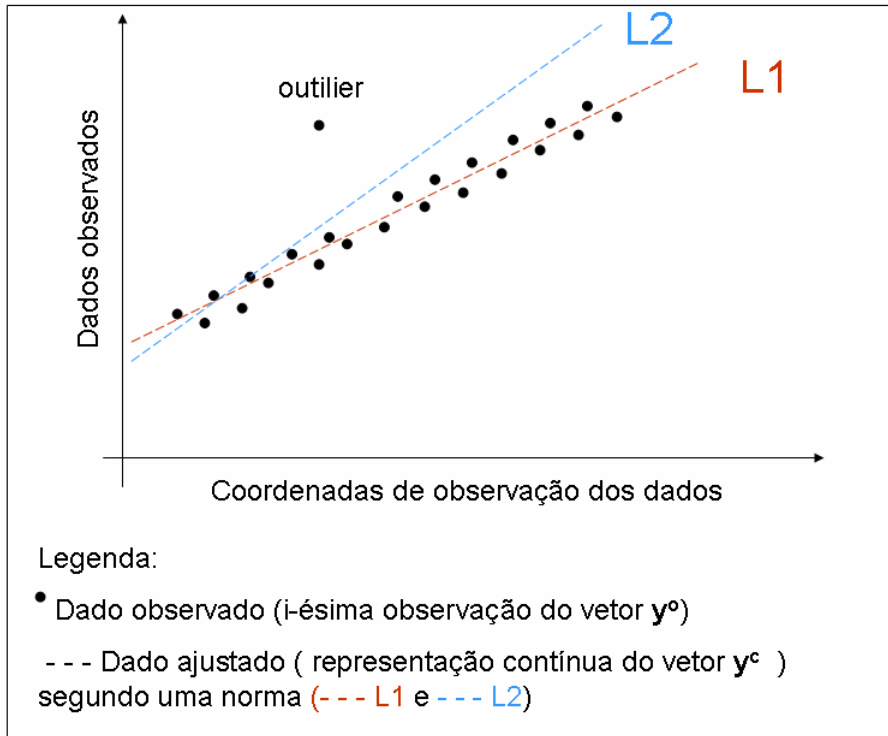


Figura 1

### Estimador da Variação Total (TV)

A função TV  $\varphi^{TV}(\mathbf{p})$  (Rudin et al., 1992) é definida como

$$\varphi^{TV}(\mathbf{p}) = \left\| \overline{\overline{\mathbf{B}}} \mathbf{p} \right\|_1, \quad (11)$$

em que  $\overline{\overline{\mathbf{B}}}$  é uma matriz representando o operador discreto de primeiras derivadas em relação as direções horizontais de distribuição dos parâmetros e  $\|\cdot\|_1$  denota norma  $\ell_1$ . Note que o produto  $\overline{\overline{\mathbf{B}}} \mathbf{p}$  representa a diferença entre parâmetros fisicamente adjacentes. Portanto, a menos de uma constante, o

produto  $\overline{\mathbf{B}} \mathbf{p}$  quantifica uma aproximação da primeira derivada<sup>1</sup> da função contínua dos parâmetros. Veja o tópico 10 para relembrar os detalhes. A equação 11 representa o regularizador de Tikhonov de ordem 1 mas usando-se a norma L1.

A função TV não penalize as discontinuidades da distribuição espacial do vetor de parâmetros  $\mathbf{p}$  de um modelo interpretativo (Vogel and Oman, 1996; 1998). Então, minimizando-se a função TV introduziremos a informação a prior que a distribuição espacial do vetor de parâmetros não será suave, mas descontínua.

Usando a definição da norma  $\ell_1$ , a função TV function dada na equação 11 pode ser reescrita como

$$\varphi^{TV}(\mathbf{p}) = \sum_{l=1}^L |p_i - p_j|, \quad (12)$$

em que  $l$  entende-se pelo  $l$ -ésimo par,  $p_i$  e  $p_j$  de parâmetros espacialmente adjacentes em relação as direções de distribuição destes parâmetros e  $L$  é o número total de pares de parâmetros espacialmente adjacentes.

Como a função TV  $\varphi^{TV}(\mathbf{p})$  não é diferenciável quando  $p_i = p_j$ , em geral usa-se a aproximação

---

<sup>1</sup>. Derivada de uma função  $f(x)$  em relação a  $x$ , considerando  $x_1$  como um número particular no domínio de  $f(x)$  é dada como:

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x_1 + \Delta x) - f(x_1)}{\Delta x}$$

$$\varphi^{TV}(\mathbf{p}) \approx \varphi_{\beta}^{TV}(\mathbf{p}) = \sum_{l=1}^L \left[ (p_i - p_j)^2 + \beta \right]^{1/2} \quad (13)$$

Proposta por Acar and Vogel (1994), em que  $\beta$  é um valor pequeno positivo. Assim a função  $\varphi_{\beta}^{TV}(\mathbf{p})$  evita dificuldades associadas com a não diferenciabilidade da função  $\varphi^{TV}(\mathbf{p})$  através da aproximação dos valores absolutos da função original  $\varphi^{TV}(\mathbf{p})$  por uma função suave que remove a discontinuidade da derivada. A função  $\varphi_{\beta}^{TV}(\mathbf{p})$  (equação 13) é em geral usada no problema vinculado não linear de minimizar

$$\lambda(\mathbf{p}) = \|\boldsymbol{\varepsilon}\|_2^2 + \mu(\delta) \varphi_{\beta}^{TV}(\mathbf{p})$$

Para detalhes veja a tese de doutorado do aluno do Observatório Nacional: Cristiano Mendel Martins (Martins, C.M. 2009)

## Referencias

- Acar, R., and C. R. Vogel, 1994, Analysis of total variation penalty methods: Inverse Problems, **10**, 1217–1229.
- Martins, C.M. 2009, Inversão gravimétrica do relevo 3d de bacias sedimentares e da variação da densidade usando informação a priori sobre o ambiente geológico: Tese de doutorado do Observatório Nacional.
- Rudin, L., S. Osher, and E. Fatemi, 1992, Nonlinear total variation based noise removal algorithms: Physica D, **60**, 259–68.
- Vogel, C. R., and M. E. Oman, 1996, Iterative methods for total variation denoising: SIAM Journal of Scientific Computing, **17**, 227–238.



\_\_\_\_\_, 1998, Fast, robust total variation-based reconstruction of noisy blurred images: IEEE Transactions on Image Processing, **7**, 813-824.