

Abordagem geométrica para imputação de dados em modelos lineares

Leandro da Silva Pereira¹

Lucas Monteiro Chaves²

Devanil Jaques Souza³

4

1 Introdução

Experimentos em que se tem perda de parcelas ou dados faltantes são bastante comuns. Nestas situações, é usual se imputar o valor destes. Este trabalho apresenta uma interpretação geométrica, em termos de subespaços vetoriais e projeções ortogonais, para um método de imputação de dados, proposto por Kruskal(1961).

2 Material e métodos

Um conjunto de dados com n observações pode ser entendido como a realização de um vetor aleatório em \mathbb{R}^n , denotado por $\mathbf{Y} = (Y_1, \dots, Y_n)'$. Suponha, sem perda de generalidade, que as primeiras m coordenadas de \mathbf{Y} sejam observáveis e as restantes $n - m$ sejam não observáveis ou faltantes. Desta forma se tem $\mathbf{Y} = \mathbf{Y}^1 + \mathbf{Y}^2$ em que \mathbf{Y}^1 é a parte observável e \mathbf{Y}^2 a não observável. Para o modelo $\mathbf{Y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}$ que leva em conta todos os dados, a matriz \mathbf{X} pode ser vista como a transformação linear $\mathbf{X} : \mathbb{R}^p \mapsto \mathbb{R}^n = \mathbb{R}^m \oplus \mathbb{R}^{n-m}$. Se $\text{Im}(\mathbf{X})$ é o subespaço definido pela imagem da transformação e P_W é a projeção ortogonal em um subespaço W , ficam definidos os subespaços $\Omega_1 = P_{\mathbb{R}^n}(\text{Im}(\mathbf{X}))$ e $\Omega_2 = P_{\mathbb{R}^{n-m}}(\text{Im}(\mathbf{X}))$. Observe que, em geral, $\text{Im}(\mathbf{X}) \neq \Omega_1 + \Omega_2$. O que se tem é que $\text{Im}(\mathbf{X}) \subset \Omega_1 + \Omega_2$. Vamos considerar a seguinte hipótese: $\dim(\text{Im}(\mathbf{X})) = \dim(\Omega_1)$. Estatisticamente, tal fato significa que o número de parâmetros do modelo, mesmo com a ocorrência de dados faltantes, não se altera. $P_{\mathbb{R}^m} : \text{Im}(\mathbf{X}) \mapsto \Omega_1$ é linear e sobrejetiva. $P_{\mathbb{R}^m} : \text{Im}(\mathbf{X}) \mapsto \Omega_1$ Como $\dim(\text{Im}(\mathbf{X})) = \dim(\Omega_1)$ então a projeção também é injetiva, isto é, $\ker(P_{\mathbb{R}^m})$ restrito à $\text{Im}(\mathbf{X})$ é o subespaço nulo. Logo, se

¹DEX-UFLA. e-mail: lspleandro_2@hotmail.com.br

²DEX - UFLA. e-mail: lucas@dex.ufla.br

³DEX - UFLA. e-mail: devaniljaques@dex.ufla.br

⁴Agradecimento à FAPEMIG pelo apoio financeiro.

$\mathbf{v}_1 + \mathbf{v}_2 \in \text{Im}(\mathbf{X})$ e $\mathbf{v}_1 + \mathbf{v}'_2 \in \text{Im}(\mathbf{X})$ então

$$\begin{aligned} (\mathbf{v}_1 + \mathbf{v}_2) - (\mathbf{v}_1 + \mathbf{v}'_2) &\in \ker(\mathbf{P}_{\mathbb{R}^m}) \\ \mathbf{v}_2 - \mathbf{v}'_2 &= \mathbf{0} \\ \mathbf{v}_2 &= \mathbf{v}'_2, \end{aligned}$$

isto é, para cada vetor \mathbf{v}_1 em Ω_1 existe um único vetor \mathbf{v}_2 em Ω_2 tal que $\mathbf{v}_1 + \mathbf{v}_2 \in \text{Im}(\mathbf{X})$, e portanto o *kernel* da projeção pertence à $\text{Im}(\mathbf{X})$. Pode-se então definir uma transformação linear $\mathbf{A} : \mathbb{R}^n \mapsto \mathbb{R}^n$ tal que $\mathbf{A}(\mathbf{v}_1) = \mathbf{v}_2$ e $\mathbf{A}(\mathbf{x}) = \mathbf{0}$ se $\mathbf{x} \in \Omega_1^\perp$.

Seja a transformação linear $\mathbf{I} : \Omega_1 \mapsto \mathbb{R}^n$, inclusão, isto é, $\mathbf{I}(\mathbf{v}_1) = \mathbf{v}_1$, $\mathbf{v}_1 \in \Omega_1$. Desta forma, pode-se escrever, utilizando um certo abuso de notação, que $\text{Im}(\mathbf{X}) = (\mathbf{I} + \mathbf{A})\Omega_1$.

Se μ é o vetor de médias geral, seu estimador de Gauss-Markov é $\hat{\mu} = \mathbf{P}_{\text{Im}(\mathbf{X})}(\mathbf{Y})$. Como as m últimas componentes de \mathbf{Y} são dados faltantes, não é possível obter $\hat{\mu}$. A estratégia então é obter o estimador de Gauss-Markov para os dados observados utilizando o sub-modelo dado por Ω_1 , isto é, $\hat{\mu}^1 = \mathbf{P}_{\text{Im}(\mathbf{X})}(\mathbf{Y}^1)$, sendo $\mu^1 = \mathbf{E}[\mathbf{Y}^1]$. Para $\hat{\mu}^1$ existe um único $\hat{\mu}^2$ dado por $\hat{\mu}^2 = \mathbf{A}(\hat{\mu}^1)$. Imputando-se o valor $\hat{\mu}^2$ obtém-se os dados com imputação $\mathbf{Y}^1 + \hat{\mu}^2$. Toma-se agora como estimador de μ o vetor obtido por

$$\hat{\mu} = \mathbf{P}_{\text{Im}(\mathbf{X})}(\mathbf{Y}^1 + \hat{\mu}^2) \quad (1)$$

A representação destes vetores e suas projeções nos subespaços discutidos podem ser observados na Figura 1.

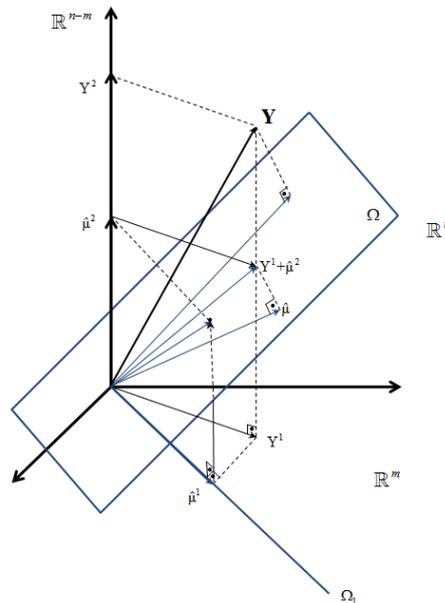


Figura 1: Projeções para dados faltantes

3 Resultados e discussão

A abordagem geométrica permite demonstrar de maneira trivial o resultado.

Proposição 1. *O processo de imputação (1) diminui a soma de quadrado dos erros.*

Demonstração. No subespaço $\text{Im}(X)$, o ponto mais próximo de $(\mathbf{Y}^1 + \hat{\mu}^2)$ é dado pela projeção ortogonal $\hat{\mu} = P_{\text{Im}(X)}(\mathbf{Y}^1 + \hat{\mu}^2)$. Note que $\hat{\mu}^1 + \hat{\mu}^2 \in \text{Im}(X)$.

Logo,

$$\left\| \underbrace{(\mathbf{Y}^1 + \hat{\mu}^2)}_{\text{dados com imputação}} - \underbrace{(\hat{\mu}^2 + \hat{\mu}^1)}_{\in \text{Im}(X)} \right\| = \|\mathbf{Y}^1 - \hat{\mu}^1\|. \text{ Assim, qualquer ponto que pertença à } \text{Im}(X)$$

e que seja diferente de $\hat{\mu}$ terá uma distância maior até o vetor $(\mathbf{Y}^1 + \hat{\mu}^2)$, o que implica em

$$\|(\mathbf{Y}^1 + \hat{\mu}^2) - \hat{\mu} = P_{\text{Im}(X)}(\mathbf{Y}^1 + \hat{\mu}^2)\| \leq \|\mathbf{Y}^1 - \hat{\mu}^1\|.$$

□

Exemplo de aplicação: delineamento inteiramente casualizado (DIC) com 3 tratamentos e 4 repetições $y_{ij} = \mu + \tau_i + \varepsilon_{ij}$, com $i = \{1, 2, 3\}$ e $j = \{1, 2, 3, 4\}$. O vetor de dados

$$\mathbf{y} = (y_{11}, y_{12}, \dots, y_{14}, y_{21}, \dots, y_{24}, y_{31}, \dots, y_{34})'.$$

Suponha a última observação perdida. O vetor de dados observado

$$\mathbf{y}^1 = (y_{11}, y_{12}, \dots, y_{14}, y_{21}, \dots, y_{24}, y_{31}, \dots, y_{33}, 0)'.$$

Tem-se

$$\text{Im}(X) = \{(a, a, a, a, b, b, b, b, c, c, c, c)'\}, a, b, c \in \mathbb{R},$$

$$\Omega_1 = \{(a, a, a, a, b, b, b, b, c, c, c, 0)'\}, a, b, c \in \mathbb{R},$$

$$\Omega_2 = \{(0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, d)'\}, d \in \mathbb{R}.$$

A transformação linear natural para este problema é $\mathbf{A} : \Omega_1 \mapsto \Omega_2$.

$$\mathbf{A}((a, a, a, a, b, b, b, b, c, c, c, 0)') = (0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, c)'$$

$$\bar{y}_1 = \frac{y_{11} + \dots + y_{14}}{4}, \bar{y}_2 = \frac{y_{21} + \dots + y_{24}}{4} \text{ e } \bar{y}_3 = \frac{y_{31} + \dots + y_{33}}{3},$$

$$P_{\Omega_1}(\mathbf{y}^1) = \left(\underbrace{\bar{y}_1, \dots, \bar{y}_1}_4, \underbrace{\bar{y}_2, \dots, \bar{y}_2}_4, \underbrace{\bar{y}_3, \dots, \bar{y}_3}_3, 0 \right)' = \hat{\mu}^1.$$

$$\mathbf{A}(\hat{\mu}^1) = (0, 0, 0, 0, 0, 0, 0, 0, 0, 0, \bar{y}_3)' = \hat{\mu}^2.$$

O vetor de dados com imputação

$$\begin{aligned} \mathbf{y}^* &= \hat{\mu}^2 + \mathbf{y}^1 \\ &= (0, 0, \dots, 0, \bar{y}_3)' + (y_{11}, \dots, y_{14}, y_{21}, \dots, y_{24}, y_{31}, \dots, y_{33}, 0)' \\ &= (y_{11}, \dots, y_{14}, y_{21}, \dots, y_{24}, y_{31}, \dots, y_{33}, \bar{y}_3)' . \end{aligned}$$

Portanto, $P_{\text{Im}(X)} \mathbf{y}^* = \left(\overbrace{\bar{y}_1, \dots, \bar{y}_1}^4, \overbrace{\bar{y}_2, \dots, \bar{y}_2}^4, \overbrace{\alpha, \dots, \alpha}^4 \right)' = \hat{\mu}$, em que

$$\begin{aligned} \alpha &= \frac{y_{31} + y_{32} + y_{33} + \frac{y_{31} + y_{32} + y_{33}}{3}}{4} \\ &= \frac{\frac{3(y_{31} + y_{32} + y_{33}) + (y_{31} + y_{32} + y_{33})}{3}}{4} \\ &= \frac{y_{31} + y_{32} + y_{33}}{3} . \end{aligned}$$

4 Conclusões

O método geométrico para abordagem de dados faltantes em modelos lineares é intuitivo e natural, permitindo uma manipulação eficiente de soma de quadrados.

Referências

- [1] SAVILLE, D. J; WOOD, G. L. **Statistical Methods: The geometric approach**. New York: Springer-Verlag, 1991. 560 p.
- [2] KRUSKAL, W. **The coordinate-free approach to Gauss-Markov estimation and its application to missing and extra observations.** Proc. Fourth Berkeley Symp. Math. Statist. Probab., v.1, p.435-451, 1961.
- [3] SILVEIRA, F. G; COSTA, L. A; PEREIRA, L. S.; et. al.; **Uma demonstração geométrica para uma identidade de Fisher para o modelo de dois fatores**. Revista Brasileira de Biometria. v. 30, n. 2, p. 199-222, 2012.