

CAP 4 - Análise de Resíduos e Diagnósticos

①

No modelo $\hat{y} = \mathbf{x}\theta + \varepsilon$, os elementos ε_i do vetor ε são as diferenças $\varepsilon_i = y_i - \hat{y}_i$ e representam a variabilidade natural dos dados. Assumimos que os erros ε_i são independentes e $\varepsilon_i \sim N(0, \sigma^2)$.

Se as suposições são violadas, temos:

- falhas sistemáticas: não linearidade; heterocedasticidade; dependência dos resíduos; não normalidade, etc
- falhas isoladas: presença de pontos atípicos

Ajustado um modelo a um conjunto de dados, o material básico que será utilizado para verificar as pressuposições, são:

- valores estimados \hat{y}_i ;
- os resíduos;
- estimativa da variância residual $\hat{\sigma}^2 = Qm_{Res}$
- elementos da diagonal da matriz H (Hat matrix)
$$H = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \quad h_{ii} \Rightarrow leverage$$

Em relação às técnicas usadas para verificação do ajuste do modelo, estas podem ser:

- formais: aninhar um modelo sob pesquisa em uma classe maior de modelos \rightarrow trabalhar a ideia de delegação de pontos
- informais: inspeção de gráficos.

Tipos de Resíduos

① Resíduos Ordinários (r_i ou $\hat{\epsilon}_i$)

Dado o ajuste do modelo $\hat{y} = X\theta + \xi$ e $\hat{Y} = HY$, o resíduo ordinário é dado por: $\hat{\xi}_i = y_i - \hat{y}_i = (I - H)y_i$.

Enquanto que os erros $\xi \sim N(\phi, I\sigma^2)$ $\begin{cases} E(\xi_i) = 0 \\ \text{cov}(\xi_i, \xi_j) = 0 \\ \text{var}(\xi_i) = \sigma^2 \end{cases}$

isso não ocorre com $\hat{\xi}_i$, ou seja,

$$\text{Var}(\hat{\xi}) = \text{Var}[(I - H)y] = (I - H)\text{Var}(y)(I - H)^T = (I - H)(I - H)^T\sigma^2 = (I - H)\sigma^2$$

$$\text{Var}(\hat{\xi}) = \begin{bmatrix} 1-h_{11} & -h_{12} & \dots & -h_{1n} \\ -h_{21} & 1-h_{22} & \dots & -h_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ -h_{n1} & -h_{n2} & \dots & 1-h_{nn} \end{bmatrix} \sigma^2$$

SIMÉTRICA

$\Rightarrow \begin{cases} \text{Var}(\hat{\epsilon}_i) = (1-h_{ii})\sigma^2 \\ \text{cov}(\hat{\epsilon}_i, \hat{\epsilon}_j) = -h_{ij}\sigma^2 \\ \text{cov}(\hat{\epsilon}_i, \hat{\epsilon}_i) = \text{Var}(\hat{\epsilon}_i) \end{cases}$

② Resíduos Padronizados (r_{pi})

Por definição os r_{pi} são obtidos por:

$$r_{pi} = \frac{\hat{\epsilon}_i}{s}, \text{ em que: } s = \sqrt{Q_{\text{res}}}$$

- permitem uma avaliação de valores discrepantes pois não são esperados valores maiores (em módulo) de que 2 ou 3;
- como a média de $\hat{\epsilon}_i$ é nula e s^2 é uma estimativa aproximada de sua variância espera-se que r_{pi} tenha uma distribuição t aproximada

→ $\hat{\epsilon}_i$ e s não são independentes
 $\rightarrow \text{Var}(\hat{\epsilon}_i) = (1-h_{ii})\sigma^2 \rightarrow$ não considerado na dist t.

③ Resíduos Estudentizados Internamente (rse_i)

$$rse_i = \frac{\hat{e}_i}{\sqrt{\text{Var}(\hat{e}_i)}} = \frac{\hat{e}_i}{\sqrt{(1-h_{ii})Q_{\text{Res}}}}$$

- são mais sensíveis que os anteriores por considerarem variâncias distintas
- valores discrepantes podem afetar profundamente a variância residual
- Numerador (\hat{e}_i) e denominador (Q_{Res}) não variáveis dependentes $\Rightarrow \text{Cov}(\hat{e}_i, Q_{\text{Res}}) \neq 0$.

④ Resíduos Estudentizados Externamente ($rse(i)$)

$$rse(i) = \frac{\hat{e}_i}{s(i)\sqrt{1-h_{ii}}} \quad \text{em que: } s(i) \text{ é o desvio-padrão residual sem a } i\text{-ésima observação.}$$

\Leftrightarrow livre da influência da i -ésima obs.

Sob normalidade, $rse(i)$ tem distribuição t-Student

com $(n-p-1)$ df.

ATENÇÃO: $|rse(i)| > t(n-p-1)$ percentil da t($n-p-1$)

ESTATÍSTICAS PARA DIAGNÓSTICOS

As observações podem ser classificadas, basicamente, como:

1) OBS INCONSISTENTES \Rightarrow resíduo grande

$$|rse(i)| \geq t(n-p-1) \text{ percentil } \frac{\alpha}{2n} \quad \begin{array}{l} \text{casos (a) e (c)} \\ \text{GRÁFICO APOSTILA} \end{array}$$

OBS: Dentro dos inconsistentes, se o leverage (h_{ii}) for pequeno \Rightarrow observação é um OUTLIER (caso a)

2) PONTO DE ALAVANCA: observações com leverage grande.

$$\text{grande} \Rightarrow h_{ii} \geq \frac{2p}{n}$$

visualmente: h_{ii} é grande quando a observação está relacionada a um x_i distante de \bar{x} .

Os pontos de alavanca podem ser:

→ Bom (se consistente) $\rightarrow r_{SC(i)}$ pef.

→ Ruim (se inconsistente) $\rightarrow r_{SC(i)}$ gde.

3) INFLUENTE: uma observação é considerada influente quando sua omissão (no conjunto de dados) resulta em mudanças significativas em alguns aspectos do modelo

\rightarrow no ajuste geral

\rightarrow nas estimativas dos parâmetros

As estatísticas:

① LEVERAGE (h_{ii}): mede a distância de uma observação em relação às demais, em termos da variável X .

No caso de uma RLS e $x_i = x_i - \bar{x}$, temos:

$$\begin{aligned}
 H &= \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ n & x_n \end{bmatrix} \begin{bmatrix} 1/n & 0 & \cdots & 0 \\ 0 & 1/\sum x_i^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1/n \end{bmatrix} \begin{bmatrix} 1 & \cdots & 1 \\ x_1 & \cdots & x_n \\ \vdots & \vdots & \vdots \\ 1 & \cdots & 1 \end{bmatrix} = \begin{bmatrix} 1/n & x_1/\sum x_i^2 & \cdots & x_n/\sum x_i^2 \\ \vdots & \vdots & \ddots & \vdots \\ 1/n & x_n/\sum x_i^2 & \cdots & x_1/\sum x_i^2 \end{bmatrix} \\
 &= \begin{bmatrix} \frac{1}{n} + \frac{x_1^2}{\sum x_i^2} & \frac{1}{n} + \frac{x_1 x_2}{\sum x_i^2} & \cdots & \frac{1}{n} + \frac{x_1 x_n}{\sum x_i^2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{n} + \frac{x_n^2}{\sum x_i^2} & \cdots & \frac{1}{n} + \frac{x_2 x_n}{\sum x_i^2} & \end{bmatrix}
 \end{aligned}$$

Portanto,

$$h_{ii} = \frac{1}{n} + \frac{x_i^2}{\sum x_i^2} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum x_i^2}$$

x_i longe $\bar{x} \Rightarrow h_{ii} \uparrow$

e

$$h_{ii} = \frac{1}{n} + \frac{x_i x_{ii}}{\sum x_i^2} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_{ii} - \bar{x})}{\sum x_i^2}$$

$x_i = \bar{x} \Rightarrow h_{ii} = \frac{1}{n}$

↳ qdb tem β_0

Para um modelo com intercepto: $\frac{1}{n} \leq h_{ii} \leq 1$

Para um modelo passa pela origem (sem β_0): $0 \leq h_{ii} \leq 1$

Além disso, $\text{tr}(H) = \sum_{i=1}^n h_{ii} = p \Rightarrow \bar{h} = p/n$ (valor médio)
 ↳ são consideradas observações influentes aquelas cujo h_{ii} supera em duas ou três vezes a sua média.

Lembrando que,

$$\hat{y} = Hy = \begin{bmatrix} h_{11} & h_{12} & \dots & h_{1n} \\ \vdots & h_{22} & & \vdots \\ \vdots & & \ddots & \vdots \\ h_{nn} & \dots & \dots & h_{nn} \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$\hat{y}_1 = h_{11} y_1 + h_{12} y_2 + \dots + h_{1n} y_n$$

↳ \hat{y}_1 estimado é uma média ponderada de todas as observações $\Rightarrow h_{11}$ são os pesos de ponderação

↳ observações com $h_{ii} \geq 2p/n$ → merecem atenção
 ↳ NO R: $h_{ii} \geq 3p/n$

② DFBeta: mede a alteração no vetor $\hat{\theta}$ ao se retirar a i -ésima observação na análise.

$$\text{DFBeta}(i) = \hat{\theta} - \hat{\theta}_{(i)}$$

NO R: influente se $|\text{DFBeta}(i)| > 1$

③ DFFITS: mede a alteração provocada no valor ajustado \hat{y}_i ao retirar a i -ésima observação da análise.

$$DFFITS = \frac{\hat{y}_i - \hat{y}_{(i)}}{\sqrt{h_{ii} s^2_{(i)}}} = \left(\frac{h_{ii}}{1-h_{ii}} \right)^{1/2} r_{Si(i)}$$

↳ potencial de influência

Merecem atenção as obs com $|DFFITS| > 2\sqrt{\frac{p}{n}}$

$$\text{No R: } |DFFITS| > 3\sqrt{\frac{p}{n-p}}$$

④ DISTANCIA DE COOK: também mede o afastamento do vetor de estimativas provocado pela retirada da i -ésima observação

$$D(i) = \frac{h_{ii}}{(1-h_{ii})} \cdot \frac{1}{p} \cdot (r_{Si})^2 \quad \text{No R: } D(i) > F_{0.05; p, n-p}$$

⑤ COV-RATIO: mede a influencia da i -ésima observação na var($\hat{\theta}$)

$$\text{COV-RATIO} = \frac{s^2_{(i)} \det |(X'_{(i)} X_{(i)})^{-1}|}{s^2 \det |(X' X)^{-1}|}$$

$$\text{No R: } |\text{COV-RATIO}| > \frac{3p}{n-p} \quad \text{ou } \frac{3p}{n}$$

GRÁFICOS

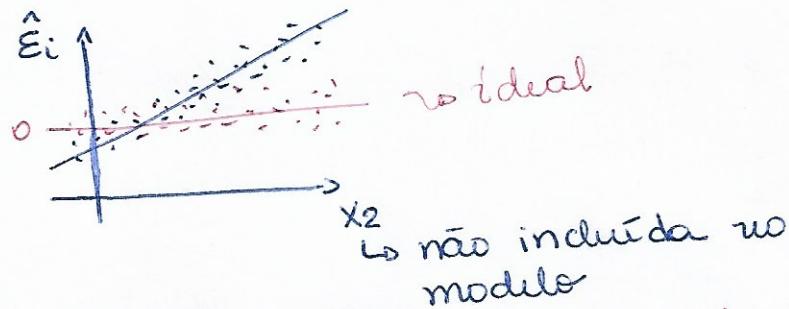
① RESÍDUOS VERSUS X_{fora}

Relações entre os resíduos do modelo ajustado e uma variável ainda não incluída no modelo.

Exemplo: $y \rightsquigarrow$ variável resposta

x_1 e $x_2 \rightsquigarrow$ variáveis regressoras

$$y_i = \beta_0 + \beta_1 x_{1i} + \epsilon_i \Rightarrow \hat{\epsilon}_i = y_i - \hat{y}_i$$



ALTERNATIVA MELHOR: gráfico da variável adicionada

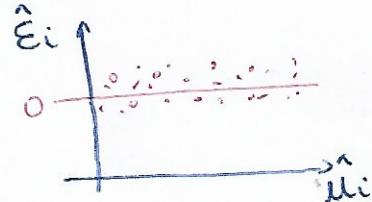
Exemplo: $y \rightsquigarrow$ variável resposta

x_1 e $x_2 \rightsquigarrow$ variáveis regressoras

$$y_i = \beta_0 + \beta_1 x_{1i} + \epsilon_i \Rightarrow \hat{\epsilon}_i \text{ (resíduo)}$$

$$x_{2i} = \beta_0 + \beta_1 x_{1i} + u_i \Rightarrow \hat{u}_i \text{ (resíduo)}$$

$$\hat{\epsilon}_i \sim \hat{u}_i$$



coef. angular não significativo ($\beta_1 = 0$) de uma reta que passa pela origem ($\beta_0 = 0$)

2) RESÍDUOS VERSUS X_{dentro}

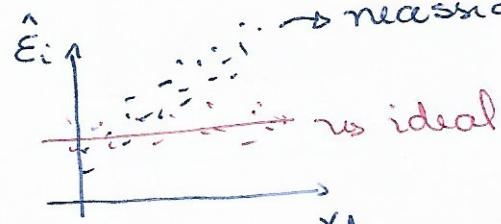
Verificar se ainda existe uma relação sistemática entre os resíduos de um modelo e a X_j já incluída no modelo \Rightarrow inclusão de termos de maior ordem (x^2)

EXEMPLO: $y \rightsquigarrow$ variável resposta

x_1 e $x_2 \rightsquigarrow$ variáveis regressoras

$$y_i = \beta_0 + \beta_1 x_{1i} + \epsilon_i \Rightarrow \hat{\epsilon}_i$$

\rightarrow necessidade de maior ordem



\rightarrow já incluída no modelo

ALTERNATIVA MELHOR: gráfico de resíduo + componente

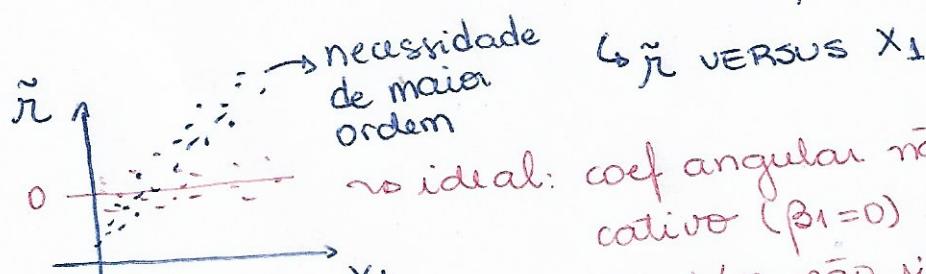
EXEMPLO: $y \rightsquigarrow$ variável resposta

x_1 e $x_2 \rightsquigarrow$ variáveis regressoras

$$y_i = \beta_0 + \beta_1 x_{1i} + \epsilon_i \Rightarrow \hat{\epsilon}_i$$

\curvearrowleft adiciona o componente $\hat{\beta}_1 x_{1i}$

$$\tilde{\epsilon}_i = \hat{\epsilon}_i + \hat{\beta}_1 x_{1i}$$



\rightarrow necessidade de maior ordem

\rightarrow ideal:

$\hookrightarrow \tilde{\epsilon}_i$ VERSUS X_1

coef angular não significativo ($\hat{\beta}_1=0$) e intercepto também não significativo ($\hat{\beta}_0=0$)