

ANÁLISE DE REGRESSÃO

Ralph dos Santos Silva

Instituto de Matemática
Centro de Ciências Matemáticas e da Natureza
Universidade Federal do Rio de Janeiro

Referências (livros)

- *A Modern Approach to Regression with R*, Sheather.
- *An R Companion to Applied Regression*, 3^a edição, Fox e Weisberg.
- ***Applied Regression Analysis*, 3^a edição, Draper e Smith.**
- *Applied Regression Modeling*, 3^a edição, Pardoe.
- *Bayesian and Frequentist Regression Methods*, Wakefield.
- *Introduction to Linear Regression Analysis*, 6^a edição, Montgomery, Peck e Vining.
- *Linear Regression and Correlation: A Beginner's Guide*, Hartshorn.
- *Regression Analysis and Linear Models: Concepts, Applications, and Implementation*, Darlington e Hayes.
- *Regression Analysis by Example*, 5^a edição, Chatterjee e Hadi.
- *Regression Analysis with Python*, Massaron e Boschetti.
- *Regression Analysis with R*, Ciaburro.
- *Regression & Linear Modeling: Best Practices and Modern Methods*, Osborne.
- *Regression Modeling with Actuarial and Financial Applications*, Frees.

Regressão

Em **análise de regressão** estamos interessados em estudar, por exemplo, relações da forma:

$$y = f(\mathbf{x}) + \varepsilon,$$

em que

- \mathbf{x} é um vetor de variáveis chamadas de *preditoras*, *explicativas* ou *regressoras*;
- y é a variável chamada de *dependente* ou *resposta*;
- a forma funcional $f(\mathbf{x})$ depende de quantidades desconhecidas β (*parâmetros*); e
- ε é um ruído aleatório.

Nota: evitaremos a utilização do termo **variáveis independentes** para \mathbf{x} .

Em geral, a forma funcional de $f(\mathbf{x})$ é, por hipótese, conhecida.

Caso contrário, utilizamos polinômios em \mathbf{x} para aproximar a verdadeira função $f(\mathbf{x})$.

Exemplo 1

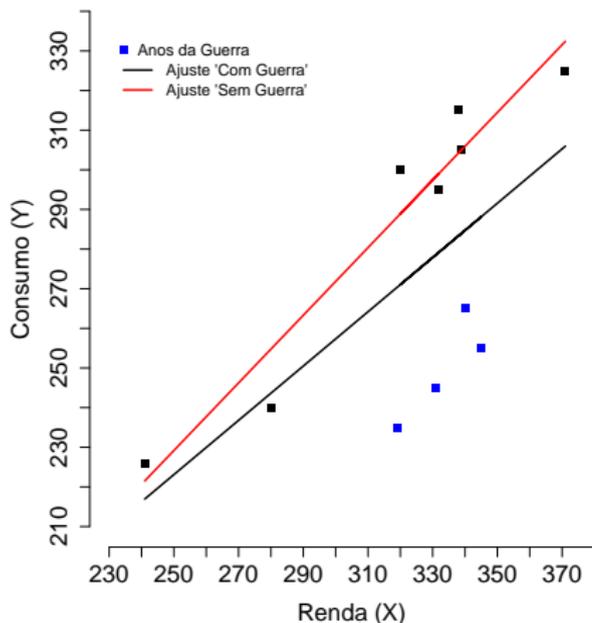


Figura 1: Relação entre a temperatura média em Fahrenheit (x) e libras de vapor por mês (y).

Ver arquivo exemplo_01.r

Regressão linear

O termo linear se refere aos parâmetros do modelo (os β 's). Por exemplo,

$$y = \beta_0 + \beta_1 x + \varepsilon;$$

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon; \text{ e}$$

$$y = \beta_0 + \beta_1 \exp(x_1) + \beta_2 \log(x_2^2) + \varepsilon,$$

são exemplos de modelos de regressão linear.

Em geral, tratamos ε e y como variáveis aleatórias enquanto \mathbf{x} é dado, isto é, a análise é condicional ao conhecimento de \mathbf{x} (regressoras).

Por hipótese, temos que **a média de ε é zero**.

Começaremos com o modelo de **regressão linear simples** dado por

$$y = \beta_0 + \beta_1 x + \varepsilon.$$

O objetivo é obter a “melhor” equação da reta que descreve a relação entre x e y .

A ideia de mínimos quadrados

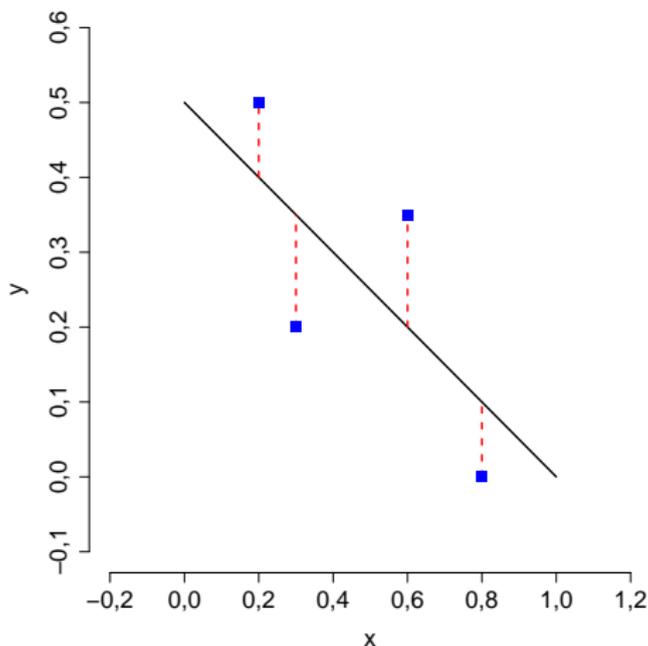


Figura 2: Ilustração da ideia de mínimos quadrados.

A **equação estimada** é denotada por $\hat{y} = b_0 + b_1x = \hat{\beta}_0 + \hat{\beta}_1x$, em que $b_0 = \hat{\beta}_0$ e $b_1 = \hat{\beta}_1$ são **estimativas pontuais** de β_0 e β_1 , respectivamente.

Mínimos quadrados

Suponha que uma amostra aleatória $((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n))$ de tamanho n seja obtida tal que

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n.$$

Definimos a função **soma de quadrados** como

$$S(\beta_0, \beta_1) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2.$$

Queremos minimizar esta soma de quadrados. Logo,

$$\begin{aligned} \frac{\partial S(\beta_0, \beta_1)}{\partial \beta_0} &= -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i), \\ \frac{\partial S(\beta_0, \beta_1)}{\partial \beta_1} &= -2 \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i). \end{aligned}$$

Agora, precisamos resolver o sistema de equações dado por

$$\begin{aligned} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) &= 0, \\ \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i) &= 0. \end{aligned}$$

Assim, temos que

$$\begin{aligned}\sum_{i=1}^n y_i - nb_0 - b_1 \sum_{i=1}^n x_i &= 0, \\ \sum_{i=1}^n x_i y_i - b_0 \sum_{i=1}^n x_i - b_1 \sum_{i=1}^n x_i^2 &= 0,\end{aligned}$$

ou

$$\begin{aligned}nb_0 + b_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i, \\ b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i.\end{aligned}$$

Ambas as representações são chamadas de **equações normais** (perpendicular ou ortogonal).

A solução do sistema acima é dado por

$$b_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{[\sum_{i=1}^n x_i][\sum_{i=1}^n y_i]}{n}}{\sum_{i=1}^n x_i^2 - \frac{[\sum_{i=1}^n x_i]^2}{n}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad e$$
$$b_0 = \bar{y} - b_1 \bar{x},$$

em que

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad e$$

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) &= \sum_{i=1}^n x_i y_i - \bar{x} \sum_{i=1}^n y_i - \bar{y} \sum_{i=1}^n x_i + n \bar{x} \bar{y} \\ &= \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} \\ &= \sum_{i=1}^n x_i y_i - \frac{[\sum_{i=1}^n x_i][\sum_{i=1}^n y_i]}{n}. \end{aligned}$$

Utilizaremos as seguintes definições:

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x})y_i = \sum_{i=1}^n x_i(y_i - \bar{y})$$

$$= \sum_{i=1}^n x_i y_i - \frac{\left[\sum_{i=1}^n x_i \right] \left[\sum_{i=1}^n y_i \right]}{n} = \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y},$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i - \bar{x})x_i$$

$$= \sum_{i=1}^n x_i^2 - \frac{\left[\sum_{i=1}^n x_i \right]^2}{n} = \sum_{i=1}^n x_i^2 - n\bar{x}^2, \quad \text{e}$$

$$S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \bar{y})y_i$$

$$= \sum_{i=1}^n y_i^2 - \frac{\left[\sum_{i=1}^n y_i \right]^2}{n} = \sum_{i=1}^n y_i^2 - n\bar{y}^2.$$

Assim, temos que $b_1 = \frac{S_{xy}}{S_{xx}}$ e $b_0 = \bar{y} - b_1\bar{x}$.

Então, agora é possível escrever a **equação ajustada** ou **predita** como

$$\hat{y} = b_0 + b_1 x.$$

Agora, substituindo $b_0 = \bar{y} - b_1 \bar{x}$ na equação acima temos que

$$\hat{y} = \bar{y} + b_1(x - \bar{x}).$$

Observe que $x = \bar{x} \Rightarrow \hat{y} = \bar{y}$.

Os dados de vapor

i	y_i	x_i
1	10,98	35,30
2	11,13	29,70
3	12,51	30,80
4	8,40	58,80
5	9,27	61,40
6	8,73	71,30
7	6,36	74,40
8	8,50	76,70
9	7,82	70,70
10	9,14	57,50
11	8,24	46,40
12	12,19	28,90
13	11,88	28,10
14	9,57	39,10
15	10,94	46,80
16	9,58	48,50
17	10,09	59,30
18	8,11	70,00
19	6,83	70,00
20	8,88	74,50
21	7,68	72,10
22	8,47	58,10
23	8,86	44,60
24	10,36	33,40
25	11,08	28,60

Cálculos para os dados de vapor

$$n = 25$$

$$\sum_{i=1}^{25} y_i = 10,98 + 11,13 + \dots + 11,08 = 235,60$$

$$\bar{y} = \frac{235,60}{25} = 9,424$$

$$\sum_{i=1}^{25} x_i = 35,3 + 29,7 + \dots + 28,6 = 1.315$$

$$\bar{x} = \frac{1.315}{25} = 52,60$$

$$\begin{aligned} \sum_{i=1}^{25} x_i y_i &= (10,98)(35,3) + (11,13)(29,7) + \dots + (11,08)(28,6) \\ &= 11.821,432 \end{aligned}$$

$$\sum_{i=1}^{25} x_i^2 = (35,3)^2 + (29,7)^2 + \dots + (28,6)^2 = 76.323,42$$

$$b_1 = \frac{11.821,432 - \frac{(1315)(235,6)}{25}}{76.323,42 - \frac{(1315)^2}{25}} = \frac{-571,128}{7.154,42} = -0,079829$$

A equação ajustada é dada por

$$\begin{aligned}\hat{y} &= \bar{y} + b_1(x - \bar{x}) = 9,4240 - 0,079829(x - 52,60) \\ &= 13,623 - 0,079829x.\end{aligned}$$

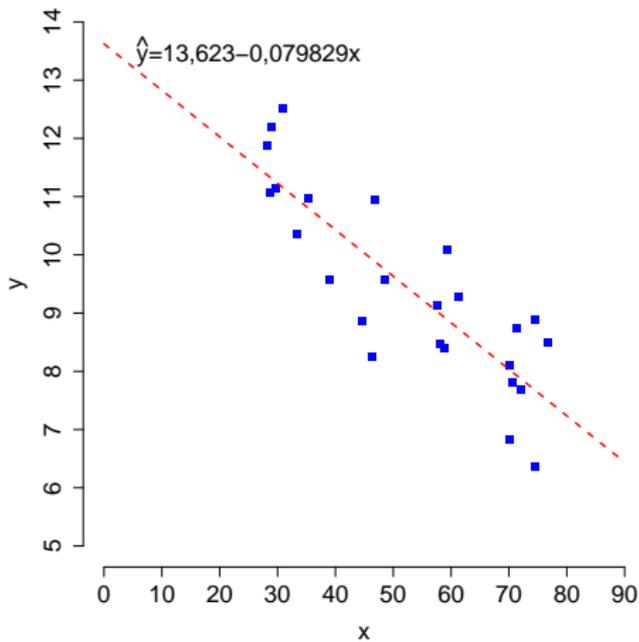


Figura 3: Equação da reta ajustada.

Resíduos

Definimos o **resíduo**, para cada observação, como diferença entre o valor observado y_i e o valor ajustado \hat{y}_i , isto é,

$$e_i = y_i - \hat{y}_i \quad \text{ou} \quad \hat{\varepsilon}_i = y_i - \hat{y}_i.$$

Importante: note que a amostra $(\hat{\varepsilon}_1, \hat{\varepsilon}_2, \dots, \hat{\varepsilon}_n)$ deve ter comportamento **similar** a uma amostra da distribuição de ε .

Agora, note que $y_i - \hat{y}_i = (y_i - \bar{y}) - b_1(x_i - \bar{x})$.

Portanto,

$$\sum_{i=1}^n \hat{\varepsilon}_i = \sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n (y_i - \bar{y}) - b_1 \sum_{i=1}^n (x_i - \bar{x}) = 0.$$

Observações, valores ajustados e resíduos

i	y_i	\hat{y}_i	$y_i - \hat{y}_i$
1	10,98	10,81	0,17
2	11,13	11,25	-0,12
3	12,51	11,16	1,35
4	8,40	8,93	-0,53
5	9,27	8,72	0,55
6	8,73	7,93	0,80
7	6,36	7,68	-1,32
8	8,50	7,50	1,00
9	7,82	7,98	-0,16
10	9,14	9,03	0,11
11	8,24	9,92	-1,68
12	12,19	11,32	0,87
13	11,88	11,38	0,50
14	9,57	10,50	-0,93
15	10,94	9,89	1,05
16	9,58	9,75	-0,17
17	10,09	8,89	1,20
18	8,11	8,03	0,08
19	6,83	8,03	-1,20
20	8,88	7,68	1,20
21	7,68	7,87	-0,19
22	8,47	8,98	-0,51
23	8,86	10,06	-1,20
24	10,36	10,96	-0,60
25	11,08	11,34	-0,26

Regressão sem a constante (sem intercepto)

Suponha que $\beta_0 = 0$, isto é, a reta passa por $(x, y) = (0, 0)$.

Então, a equação a ser estimada é dada por $y_i = \beta_1 x_i + \varepsilon_i$.

Derivando-se $S(\beta_1)$ em relação a β_1 e igualando-se a zero, temos que

$$\sum_{i=1}^n x_i (y_i - \beta_1 x_i) = 0 \quad \Rightarrow \quad b_1 = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}.$$

A equação da reta estimada é dada por $\hat{y} = b_1 x$.

No ponto $x = \bar{x}$ temos que $\hat{y} = b_1 \bar{x}$, isto é, não resulta em \bar{y} .

Além disso, em geral,

$$\sum_{i=1}^n \hat{\varepsilon}_i = \sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n (y_i - b_1 x_i) \neq 0.$$

Se $(x, y) = (0, 0)$ for verdade, então $b_0 = 0$. Consequentemente $b_1 = \frac{\bar{y}}{\bar{x}}$ que resulta em $\sum_{i=1}^n \hat{\varepsilon}_i = 0$.

Centrando os dados

Temos que $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$.

Agora, suponha que $y_i - \bar{y} = (\beta_0 + \beta_1 \bar{x} - \bar{y}) + \beta_1 (x_i - \bar{x}) + \varepsilon_i$.

Reescrevendo, obtemos que $y_i^* = \beta_0^* + \beta_1 x_i^* + \varepsilon_i$, em que

$$\begin{aligned}y_i^* &= y_i - \bar{y} \\ \beta_0^* &= \beta_0 + \beta_1 \bar{x} - \bar{y} \\ x_i^* &= x_i - \bar{x}.\end{aligned}$$

Note que $b_1 = \frac{\sum_{i=1}^n x_i^* y_i^*}{\sum_{i=1}^n [x_i^*]^2}$ e $b_0^* = \bar{y}^* - b_1 \bar{x}^* = 0$, pois $\bar{x}^* = \bar{y}^* = 0$.

Como isto sempre acontece ($b_0^* = 0$), o modelo a ser ajustado é dado por

$$y_i - \bar{y} = b_1 (x_i - \bar{x})$$

Perdemos um parâmetro (β_0). Contudo, as quantidades $(y_i - \bar{y})$, para $i = 1, 2, \dots, n$, representam somente $(n - 1)$ pedaços de informações.

Análise de variância

Temos que avaliar a variação nos dados explicada pela reta de regressão.

Considere que $y_i - \hat{y}_i = (y_i - \bar{y}) - (\hat{y}_i - \bar{y})$.

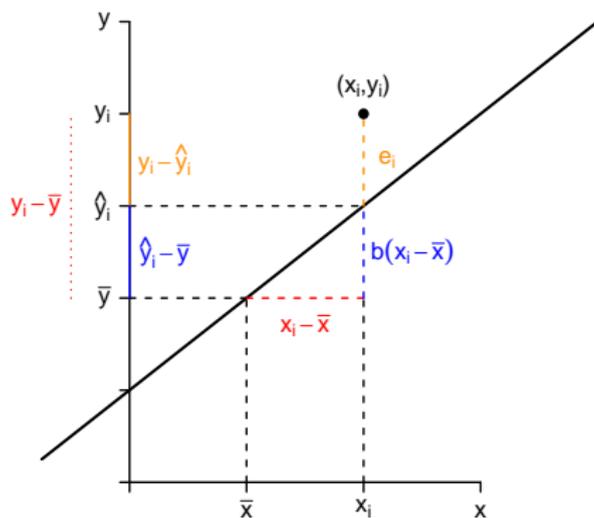


Figura 4: Decomposição de y .

Note que

$$\bar{\hat{y}} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n (b_0 + b_1 x_i) = \frac{1}{n} (nb_0 + b_1 n\bar{x}) = b_0 + b_1 \bar{x} = \bar{y}.$$

Isto implica novamente que

$$\sum_{i=1}^n e_i = \sum_{i=1}^n \hat{\varepsilon}_i = \sum_{i=1}^n (y_i - \hat{y}_i) = n\bar{y} - n\bar{y} = 0.$$

Podemos reescrever a decomposição como

$$(y_i - \bar{y}) = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i).$$

Portanto, a soma de quadrados resulta em

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

pois a soma do produto cruzado (SPC) é igual a zero, isto é,

$$\text{SPC} = 2 \sum_{i=1}^n (\hat{y}_i - \bar{y})(y_i - \hat{y}_i) = 0.$$

Agora, note que **(ver Figura 4)**

$$\hat{y}_i - \bar{y} = b_1(x_i - \bar{x}) \quad \text{e} \quad y_i - \hat{y}_i = y_i - \bar{y} - b_1(x_i - \bar{x}).$$

Assim,

$$\text{SPC} = 2 \sum_{i=1}^n b_1(x_i - \bar{x})[(y_i - \bar{y}) - b_1(x_i - \bar{x})] = 2b_1(S_{xy} - b_1 S_{xx}) = 0,$$

pois

$$b_1 = \frac{S_{xy}}{S_{xx}}.$$

Temos também que

$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n b_1^2 (x_i - \bar{x})^2 = b_1^2 S_{xx} = \frac{S_{xy}^2}{S_{xx}}.$$

A soma de quadrados é dada por

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{SQTot} = \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{SQReg} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{SQRes},$$

em que

- $SQTot$ é a soma de quadrados total;
- $SQReg$ é a soma de quadrados da regressão; e
- $SQRes$ é a soma de quadrados dos resíduos.

Tabela 1: Análise de variância.

Soma de Quadrados	Fonte	Graus de Liberdade	Média Quadrática
Regressão	$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	1	$MQReg$
Resíduo	$\sum_{i=1}^n (y_i - \hat{y}_i)^2$	$n - 2$	s^2
Total	$\sum_{i=1}^n (y_i - \bar{y})^2$	$n - 1$	s_y^2

Por definição, temos que $MQReg = \frac{SQReg}{1}$, $s^2 = \frac{SQRes}{n - 2}$ e $s_y^2 = \frac{S_{yy}}{(n - 1)}$.

Coeficiente de determinação

$$R^2 = \frac{SQReg}{SQTot} = 1 - \frac{SQRes}{SQTot} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

O R^2 mede a proporção da variação total que é explicada pela regressão.

Exemplo 1 (continuação)

Tabela 2: Análise de variância - dados de vapor.

Soma de Quadrados	Fonte	Graus de Liberdade	Média Quadrática	R^2
Regressão	45,59	1	45,59	0,714
Resíduo	18,22	23	0,792	—
Total	63,82	24	2,659	—

Intervalos de confiança e testes de hipóteses para β_0 e β_1

Seja o modelo $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, para $i = 1, 2, \dots, n$.

Agora, consideraremos as seguintes hipóteses:

HP.1: ε_i é uma variável aleatória com média 0 (zero) e variância constante e desconhecida σ_ε^2 . Temos que $E(\varepsilon_i) = 0$ e $\text{Var}(\varepsilon_i) = \sigma_\varepsilon^2$.

HP.2: ε_i e ε_j são não correlacionados para todo $i \neq j$ tal que $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$.

HP.3: $\varepsilon_i \sim \mathcal{N}(0, \sigma_\varepsilon^2)$, isto é, ε_i tem distribuição normal.

Consequentemente, temos que

$$\begin{aligned} E(y_i | x_i) &= \beta_0 + \beta_1 x_i; \\ \text{Var}(y_i | x_i) &= \sigma_\varepsilon^2; \\ \text{Cov}(y_i, y_j | x_i, x_j) &= 0, \quad \text{para todo } i \neq j; \text{ e} \\ (y_i | x_i) &\sim \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma_\varepsilon^2). \end{aligned}$$

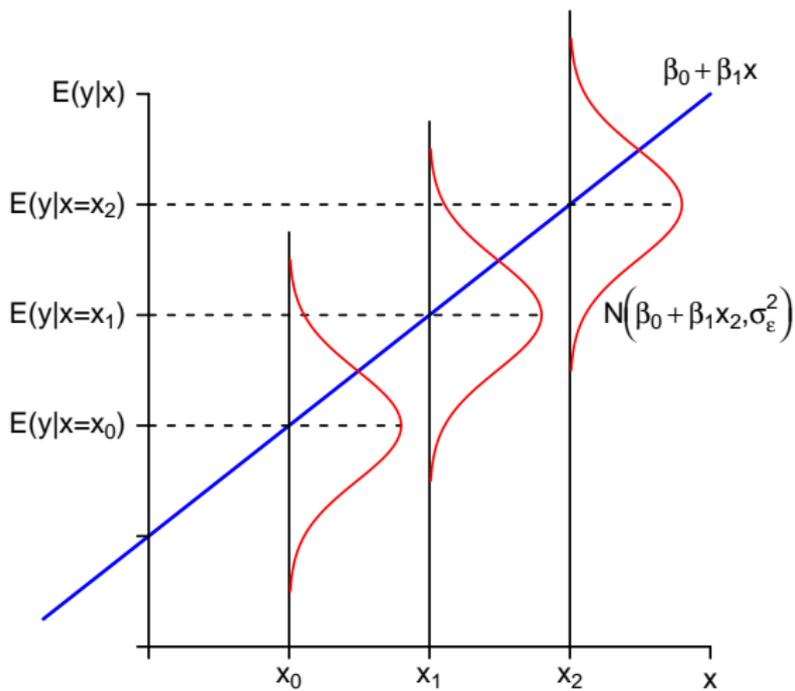


Figura 5: O modelo clássico de regressão com erros normais.

Propriedades de b_1

Temos que

$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} = \sum_{i=1}^n \omega_i y_i,$$

em que

$$\omega_i = \frac{(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Assim,

$$\begin{aligned} E(b_1 | \mathbf{x}) &= E\left(\sum_{i=1}^n \omega_i y_i \mid \mathbf{x}\right) = \sum_{i=1}^n \omega_i E(y_i | x_i) \\ &= \sum_{i=1}^n \omega_i (\beta_0 + \beta_1 x_i) = \frac{\beta_0 \sum_{i=1}^n (x_i - \bar{x}) + \beta_1 \sum_{i=1}^n (x_i - \bar{x}) x_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \beta_0 \frac{\sum_{i=1}^n (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} + \beta_1 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \beta_1. \end{aligned}$$

$$\begin{aligned} \text{Var}(b_1|\mathbf{x}) &= \text{Var}\left(\sum_{i=1}^n \omega_i y_i \mid \mathbf{x}\right) = \sum_{i=1}^n \omega_i^2 \text{Var}(y_i|x_i) \\ &= \sum_{i=1}^n \omega_i^2 \sigma_\varepsilon^2 = \sigma_\varepsilon^2 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\left[\sum_{i=1}^n (x_i - \bar{x})^2\right]^2} = \frac{\sigma_\varepsilon^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sigma_\varepsilon^2}{S_{xx}}. \end{aligned}$$

Portanto,

$$\text{DP}(b_1|\mathbf{x}) = \frac{\sigma_\varepsilon}{\left[\sum_{i=1}^n (x_i - \bar{x})^2\right]^{1/2}} = \frac{\sigma_\varepsilon}{S_{xx}^{1/2}}.$$

Assumindo que o modelo é o correto, substituímos σ_ε por s e obtemos que

$$\text{ep}(b_1|\mathbf{x}) = \widehat{\text{DP}}(b_1|\mathbf{x}) = \frac{s}{\left[\sum_{i=1}^n (x_i - \bar{x})^2\right]^{1/2}} = \frac{s}{S_{xx}^{1/2}}.$$

Se considerarmos **HP.3**, e como $b_1 = \sum_{i=1}^n \omega_i y_i$ é uma combinação linear de normais, temos que

$$(b_1|\mathbf{x}) \sim \mathcal{N}\left(\beta_1, \frac{\sigma_\varepsilon^2}{S_{xx}}\right).$$

Exemplo 2: distribuição amostral do estimador de mínimos quadrados

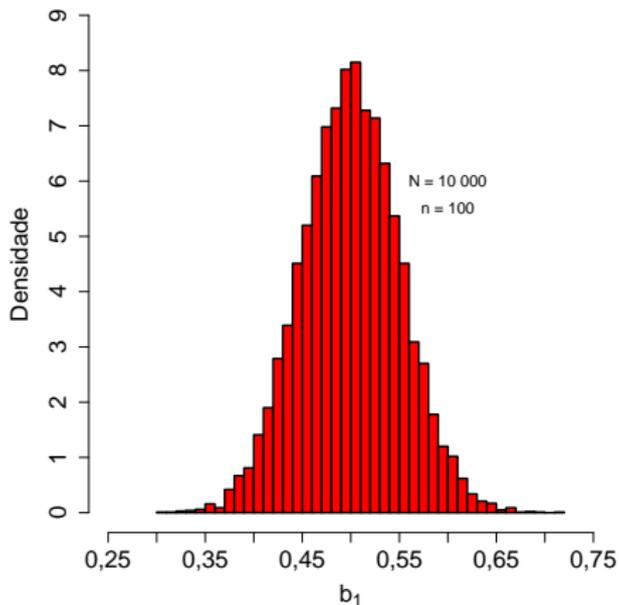
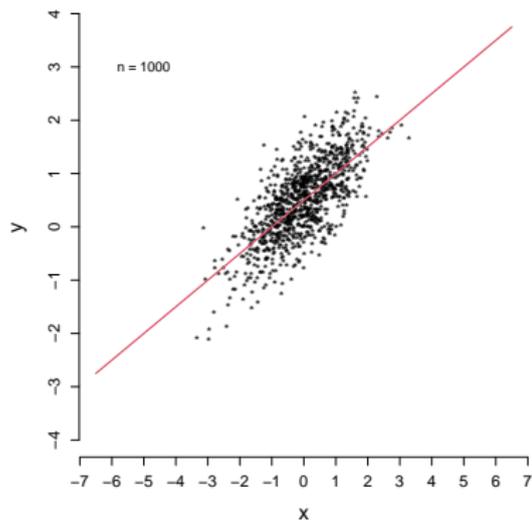


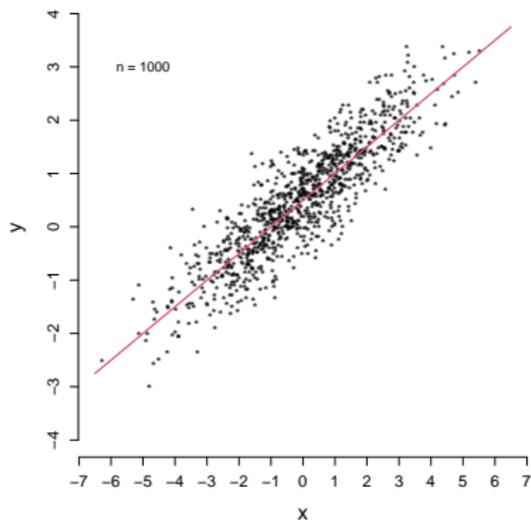
Figura 6: Distribuição amostral do estimador b_1 de β_1 na regressão $y = \beta_0 + \beta_1 x + \varepsilon$.

Ver arquivo exemplo_02.r

Exemplo 3: variância do estimador de mínimos quadrados



(a) Menor variância em x .



(b) Maior variância em x .

Figura 7: Variância do estimador b_1 de β_1 na regressão $y = \beta_0 + \beta_1 x + \varepsilon$ com $S_{xx}(1) = 1.015,76$ e $S_{xx}(2) = 3.657,90$.

Ver arquivo exemplo_03.r

Intervalo de confiança e teste de hipóteses para β_1

Se considerarmos **HP.3**, então $(b_1|\mathbf{x}) \sim \mathcal{N}\left(\beta_1, \frac{\sigma_\varepsilon^2}{S_{xx}}\right)$.

Pode ser mostrado que

$$\text{IC}_{100(1-\alpha)\%}(\beta_1) = \left(b_1 - t_{n-2, 1-\alpha/2} \frac{s}{\sqrt{S_{xx}}}, b_1 + t_{n-2, 1-\alpha/2} \frac{s}{\sqrt{S_{xx}}} \right).$$

Além disso, para o teste estatístico de hipóteses

$$\begin{cases} H_0 : \beta_1 = \beta_{1,0} \\ H_1 : \beta_1 \neq \beta_{1,0} \end{cases}$$

temos que a estatística do teste é dada por

$$t = \frac{(b_1 - \beta_{1,0})}{\text{ep}(b_1)} = \frac{(b_1 - \beta_{1,0})\sqrt{S_{xx}}}{s}.$$

Se $|t| > t_{n-2, 1-\alpha/2}$, então rejeita-se a hipótese nula.

Exemplo 1 (continuação)

O erro padrão de b_1 é dado por $ep(b_1) = \frac{s}{\sqrt{S_{xx}}} = \frac{0,8901}{\sqrt{84,58}} = 0,01052$.

Se considerarmos $\alpha = 0,05$, temos que $t_{23;0,975} = 2,069$.

Assim,

$$\begin{aligned} IC_{95\%}(\beta_1) &= (-0,0798 - 2,069 \times 0,01052; -0,0798 + 2,069 \times 0,01052) \\ &= (-0,10160; -0,05806). \end{aligned}$$

Agora, para o teste

$$\begin{cases} H_0 : \beta_1 = 0 \\ H_1 : \beta_1 \neq 0, \end{cases}$$

temos que $t = (b_1 - 0)/ep(b_1) = -0,0798/0,01052 = -7,60$.

Como $|t| = 7,60 > t_{23;0,975} = 2,069$, então rejeita-se a H_0 , isto é, β_1 é estatisticamente diferente de zero.

Concluimos que a variável x realmente ajuda explicar a variável y através da regressão linear $y = \beta_0 + \beta_1 x + \varepsilon$.

Ver arquivo [exemplo_01.r](#)

Propriedades de b_0

Temos que $b_0 = \bar{y} - b_1\bar{x}$. Assim,

$$\begin{aligned}E(b_0|\mathbf{x}) &= E(\bar{y} - b_1\bar{x}|\mathbf{x}) = E(\bar{y}|\mathbf{x}) - E(b_1|\mathbf{x})\bar{x} \\ &= \frac{1}{n}\sum_{i=1}^n E(y_i|x_i) - \beta_1\bar{x} = \frac{1}{n}\sum_{i=1}^n (\beta_0 + \beta_1 x_i) - \beta_1\bar{x} = \beta_0\end{aligned}$$

$$\begin{aligned}\text{Var}(b_0|\mathbf{x}) &= \text{Var}(\bar{y} - b_1\bar{x}|\mathbf{x}) \\ &= \text{Var}(\bar{y}|\mathbf{x}) + [\bar{x}]^2\text{Var}(b_1|\mathbf{x}) - 2\text{Cov}(\bar{y}, b_1|\mathbf{x})\bar{x} \\ &= \frac{\sigma_\varepsilon^2}{n} + [\bar{x}]^2 \frac{\sigma_\varepsilon^2}{S_{xx}} = \sigma_\varepsilon^2 \left[\frac{1}{n} + \frac{[\bar{x}]^2}{S_{xx}} \right] \\ &= \sigma_\varepsilon^2 \left[\frac{S_{xx} + n[\bar{x}]^2}{nS_{xx}} \right] = \sigma_\varepsilon^2 \left[\frac{\sum_{i=1}^n x_i^2}{nS_{xx}} \right], \quad \text{pois}\end{aligned}$$

$$\begin{aligned}\text{Cov}(\bar{y}, b_1|\mathbf{x}) &= \text{Cov}\left(\frac{1}{n}\sum_{i=1}^n y_i, \sum_{i=1}^n \omega_i y_i|\mathbf{x}\right) \\ &= \frac{1}{n}\sum_{i=1}^n \text{Cov}(y_i, \omega_i y_i|\mathbf{x}) = \frac{1}{n}\sum_{i=1}^n \omega_i \text{Var}(y_i|\mathbf{x}) \\ &= \frac{1}{n}\sum_{i=1}^n \omega_i \sigma_\varepsilon^2 = \frac{\sigma_\varepsilon^2}{n} \sum_{i=1}^n \frac{(x_i - \bar{x})}{S_{xx}} = 0.\end{aligned}$$

Portanto,

$$DP(b_0|\mathbf{x}) = \sigma_\varepsilon \left[\frac{\sum_{i=1}^n x_i^2}{nS_{xx}} \right]^{1/2}.$$

Assumindo que o modelo é o correto, substituímos σ_ε por s e obtemos que

$$ep(b_0|\mathbf{x}) = DP(\widehat{b_0}|\mathbf{x}) = s \left[\frac{\sum_{i=1}^n x_i^2}{nS_{xx}} \right]^{1/2}.$$

Se considerarmos **HP.3**, e como $b_0 = \bar{y} - b_1\bar{x}$ é uma combinação linear de normais, temos que

$$(b_0|\mathbf{x}) \sim \mathcal{N} \left(\beta_0, \sigma_\varepsilon^2 \left[\frac{\sum_{i=1}^n x_i^2}{nS_{xx}} \right] \right).$$

Intervalo de confiança e teste de hipóteses para β_0

Pode ser mostrado que

$$IC_{100(1-\alpha)\%}(\beta_0) = \left(b_0 - t_{n-2, 1-\alpha/2} s \left[\frac{\sum_{i=1}^n x_i^2}{nS_{xx}} \right]^{1/2}, b_0 + t_{n-2, 1-\alpha/2} s \left[\frac{\sum_{i=1}^n x_i^2}{nS_{xx}} \right]^{1/2} \right).$$

Além disso, para o teste estatístico de hipóteses

$$\begin{cases} H_0 : \beta_0 = \beta_{0,0} \\ H_1 : \beta_0 \neq \beta_{0,0} \end{cases}$$

temos que a estatística do teste é dada por

$$t = \frac{(b_0 - \beta_{0,0})}{\text{ep}(b_0)} = \frac{(b_0 - \beta_{0,0})\sqrt{nS_{xx}}}{s\sqrt{\sum_{i=1}^n x_i^2}}.$$

Se $|t| > t_{n-2, 1-\alpha/2}$, então rejeita-se a hipótese nula.

Tabela 3: Estimativas para os dados com **menor** dispersão em x .

Parâmetro	Estimativa	Erro padrão	Estat. t	Valor p
β_0	0,52368	0,01579	33,16	$< 10^{-15}$
β_1	0,51920	0,01566	33,15	$< 10^{-15}$

$$\hat{\sigma}_1 = 0,4992 \text{ e } R_1^2 = 0,5240.$$

Tabela 4: Estimativas para os dados com **maior** dispersão em x .

Parâmetro	Estimativa	Erro padrão	Estat. t	Valor p
β_0	0,52294	0,01580	33,09	$< 10^{-15}$
β_1	0,49683	0,00826	60,15	$< 10^{-15}$

$$\hat{\sigma}_2 = 0,4996 \text{ e } R_2^2 = 0,7838.$$

Propriedades das somas de quadrados

Suponha que $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, $i = 1, 2, \dots, n$, seja o modelo correto. Assim,

$$\begin{aligned}E(\bar{y}|\mathbf{x}) &= E\left(\frac{1}{n}\sum_{i=1}^n y_i \mid \mathbf{x}\right) = \frac{1}{n}\sum_{i=1}^n E(y_i|\mathbf{x}) \\&= \frac{1}{n}\sum_{i=1}^n (\beta_0 + \beta_1 x_i) = \beta_0 + \beta_1 \bar{x} \\E(SQ_{Tot}|\mathbf{x}) &= E\left(\sum_{i=1}^n (y_i - \bar{y})^2 \mid \mathbf{x}\right) = E\left(\sum_{i=1}^n y_i^2 - n\bar{y}^2 \mid \mathbf{x}\right) \\&= \sum_{i=1}^n E(y_i^2|x_i) - nE(\bar{y}^2|\mathbf{x}) \\&= \sum_{i=1}^n \left[\text{Var}(y_i|x_i) + E^2(y_i|x_i)\right] - n\left[\text{Var}(\bar{y}|\mathbf{x}) + E^2(\bar{y}|\mathbf{x})\right] \\&= \sum_{i=1}^n \left[\sigma_\varepsilon^2 + (\beta_0 + \beta_1 x_i)^2\right] - n\left[\frac{\sigma_\varepsilon^2}{n} + (\beta_0 + \beta_1 \bar{x})^2\right] \\&= (n-1)\sigma_\varepsilon^2 + \sum_{i=1}^n (\beta_0 + \beta_1 x_i)^2 - n(\beta_0 + \beta_1 \bar{x})^2 \\&= (n-1)\sigma_\varepsilon^2 + \beta_1^2 S_{xx}\end{aligned}$$

Continuando, temos que

$$\begin{aligned}E(SQReg|\mathbf{x}) &= E\left(\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \mid \mathbf{x}\right) = \sum_{i=1}^n E((\hat{y}_i - \bar{y})^2 \mid \mathbf{x}) \\&= \sum_{i=1}^n E(b_1^2 (x_i - \bar{x})^2 \mid \mathbf{x}) = E(b_1^2 \mid \mathbf{x}) S_{xx} \\&= \left[\text{Var}(b_1 \mid \mathbf{x}) + E^2(b_1 \mid \mathbf{x}) \right] S_{xx} = \left[\frac{\sigma_\varepsilon^2}{S_{xx}} + \beta_1^2 \right] S_{xx} \\&= \sigma_\varepsilon^2 + \beta_1^2 S_{xx}. \\E(SQRes|\mathbf{x}) &= E(SQTot - SQReg|\mathbf{x}) = E(SQTot|\mathbf{x}) - E(SQReg|\mathbf{x}) \\&= (n-1)\sigma_\varepsilon^2 + \beta_1^2 S_{xx} - \sigma_\varepsilon^2 - \beta_1^2 S_{xx} = (n-2)\sigma_\varepsilon^2.\end{aligned}$$

Consequentemente,

$$E(MQReg|\mathbf{x}) = \sigma_\varepsilon^2 + \beta_1^2 S_{xx} \quad \text{e} \quad E(s^2|\mathbf{x}) = \sigma_\varepsilon^2.$$

Se considerarmos **HP.3** e que $\beta_1 = 0$, então é possível mostrar $MQReg$ e s^2 são independentes, e que

$$\frac{MQReg}{\sigma_\varepsilon^2} \sim \chi_1^2 \quad \text{e} \quad \frac{(n-2)s^2}{\sigma_\varepsilon^2} \sim \chi_{n-2}^2.$$

Além disso, $\frac{MQReg}{s^2} \sim \mathcal{F}_{1, n-2}$.

Teste F para significância da regressão

Para o teste estatístico de hipóteses

$$\begin{cases} H_0 : \beta_1 = 0 \\ H_1 : \beta_1 \neq 0 \end{cases}$$

temos que a estatística do teste (alternativa) é dada por

$$F = \frac{MQReg}{s^2}.$$

Se $F > F_{1,n-2,1-\alpha}$, então rejeita-se a hipótese nula.

Note que os testes T e F são equivalentes (quando temos somente uma variável regressora), pois neste caso $F = T^2 \sim \mathcal{F}_{1,n-2}$.

Exemplo 1 (continuação)

Queremos testar via estatística

$$\begin{cases} H_0 : \beta_1 = 0 \\ H_1 : \beta_1 \neq 0 \end{cases}$$

Temos que

$$F = \frac{MQReg}{s^2} = \frac{45,59}{0,792} = 57,54$$

e $F_{1;23;0,95} = 4,279$. Como $F > F_{1;23;0,95}$, então rejeitamos H_0 .

Note que $t = -7,586$ e $F = t^2 = 57,54$. Então, para o modelo com uma regressora, os testes t e F geram os mesmos resultados e conclusões.

Tabela 5: Análise de variância - dados de vapor.

Soma de Quadrados	Fonte	Graus de Liberdade	Média Quadrática	F	R^2
Regressão	45,59	1	45,59	57,54	0,714
Resíduo	18,22	23	0,792	—	—
Total	63,82	24	2,659	—	—

Ver arquivo exemplo_01.r

Correlação de Pearson, b_1 e R^2

O coeficiente de correlação de Pearson entre duas variáveis aleatórias é definido por

$$\rho(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(y)\text{Var}(x)}}$$

e mede o grau de relação linear entre estas variáveis.

Dado uma amostra aleatória $((x_1, y_1), \dots, (x_n, y_n))$, temos a correlação amostral dada por

$$r(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\left[\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2 \right]^{1/2}} = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}}$$

Assim, temos que

$$b_1 = \frac{S_{xy}}{S_{xx}} = \frac{S_{xy} S_{yy}^{1/2}}{S_{xx}^{1/2} S_{xx}^{1/2} S_{yy}^{1/2}} = r(x, y) \frac{S_{yy}^{1/2}}{S_{xx}^{1/2}}$$

Então, b_1 e $r(x, y)$ têm o mesmo sinal.

Dado que $(\hat{y}_i - \bar{y}) = b_1(x_i - \bar{x})$, note que

$$\begin{aligned}R^2 &= \frac{SQReg}{SQTot} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\sum_{i=1}^n b_1^2 (x_i - \bar{x})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \\&= \frac{b_1^2 S_{xx}}{S_{yy}} = \frac{S_{xy}^2 S_{xx}}{S_{xx}^2 S_{yy}} = \left[\frac{S_{xy}}{S_{xx}^{1/2} S_{yy}^{1/2}} \right]^2 = [r(x, y)]^2; \\[r(x, y)]^2 &= \frac{\left[\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right]^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2} \\&= \frac{\left[\sum_{i=1}^n b_1 (x_i - \bar{x})(y_i - \bar{y}) \right]^2}{\sum_{i=1}^n b_1^2 (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2} \\&= \frac{\left[\sum_{i=1}^n (\hat{y}_i - \bar{y})(y_i - \bar{y}) \right]^2}{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \sum_{i=1}^n (y_i - \bar{y})^2} = [r(y, \hat{y})]^2; \quad e \\R^2 &= [r(y, \hat{y})]^2.\end{aligned}$$

(É possível testar $H_0 : \rho(x, y) = 0$ contra $H_1 : \rho(x, y) \neq 0$.)

O coeficiente de determinação R^2

Sejam \mathcal{M}_1 e \mathcal{M}_2 os seguintes modelos de regressão:

$$\mathcal{M}_1 : \quad y_i = \beta_0 + \beta_1 x_{1i} + \varepsilon_i$$

$$\mathcal{M}_2 : \quad y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i,$$

com coeficientes de determinação R_1^2 e R_2^2 , respectivamente.

Então, temos que $R_2^2 \geq R_1^2$.

O resultado diz que adicionar outra covariável ao modelo tende a melhorar a regressão (aumentar a $SQReg$) e diminuir o resíduo ($SQRes$).

Análise dos resíduos

Podemos realizar inspeções gráficas ou testes de hipóteses.

- Verificar a normalidade;
- Verificar o efeito do tempo se os dados forem ordenados no tempo;
- Verificar a homocedasticidade (variância constante);
- Verificar possíveis transformações dos dados;
- Verificar uma ordem polinomial maior do que a ajustada com os \mathbf{x} 's;
- Verificar pontos aberrantes ou de alavancagem; e
- Verificar as hipóteses atreladas ao tipo de dado.

Note que

- $\sum_{i=1}^n \hat{\varepsilon}_i = 0$;
- $\text{Cov}(x, \hat{\varepsilon}) = 0$ (por hipótese);
- Se $t = 1, 2, \dots, n$, então $\text{Cov}(t, \hat{\varepsilon}) = 0$ (por hipótese);
- $\text{Cov}(\hat{y}, \hat{\varepsilon}) = 0$, mas $\text{Cov}(y, \hat{\varepsilon}) \neq 0$; e
- $\text{Var}(y|\mathbf{x}) = \sigma^2$.

Algumas definições (pausa em regressão)

Momentos ordinários:

$$m'_\ell = E(y^\ell) = \int_{-\infty}^{\infty} y^\ell f(y) dy.$$

Definimos $\mu_y = m'_1 = E(y)$ como a média.

Momentos centrais:

$$m_\ell = E\left((y - \mu_y)^\ell\right) = \int_{-\infty}^{\infty} (y - \mu_y)^\ell f(y) dy.$$

Definimos $\sigma_y^2 = m_2 = E\left((y - \mu_y)^2\right)$ como a variância.

O coeficiente de assimetria é definido por

$$A(y) = E \left(\frac{[y - \mu_y]^3}{\sigma_y^3} \right)$$

e o coeficiente de curtose por

$$K(y) = E \left(\frac{[y - \mu_y]^4}{\sigma_y^4} \right).$$

- A quantidade $K(y) - 3$ é chamada de excesso de curtose porque $K(y) = 3$ para a distribuição normal.
- Uma distribuição com excesso de curtose positivo é dita ter caudas pesadas (leptocúrtica).
- Uma distribuição com excesso de curtose negativo é dita ter caudas leves (platicúrtica).

Momentos amostrais

Suponha que $\{y_1, y_2, \dots, y_n\}$ seja uma amostra aleatória de y com n observações.

- Média amostral: $\hat{\mu}_y = \frac{1}{n} \sum_{i=1}^n y_i$.
- Variância amostral: $\hat{\sigma}_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \hat{\mu}_y)^2$.
- O coeficiente de assimetria amostral: $\hat{A}(y) = \frac{1}{(n-1)\hat{\sigma}_y^3} \sum_{i=1}^n (y_i - \hat{\mu}_y)^3$.
- O coeficiente de curtose amostral: $\hat{K}(y) = \frac{1}{(n-1)\hat{\sigma}_y^4} \sum_{i=1}^n (y_i - \hat{\mu}_y)^4$.

- Sob a hipótese de normalidade (**hipótese de que os dados são provenientes de uma distribuição normal**), temos que

$$\widehat{A}(y) \approx \mathcal{N}(0, 6/n) \quad \text{e} \quad \widehat{K}(y) \approx \mathcal{N}(3, 24/n)$$

para n “suficientemente grande”.

- Estas aproximações para n grande podem ser utilizadas para testar a hipótese de normalidade dos dados.
- Dado uma amostra aleatória $\{y_1, y_2, \dots, y_n\}$, para testar a assimetria dos retornos, consideramos

$$H_0 : A(y) = 0$$

$$H_1 : A(y) \neq 0.$$

- A estatística da razão t -Student da assimetria amostral é

$$t = \frac{\widehat{A}(y)}{\sqrt{6/n}} \approx \mathcal{N}(0, 1).$$

- Rejeitamos H_0 ao nível α de significância se $|t| > z_{1-\alpha/2}$, em que $z_{1-\alpha/2}$ o percentil $100(1 - \alpha/2)$ da distribuição normal padrão.

- Para testar o excesso de curtose dos retornos, consideramos

$$H_0 : K(y) - 3 = 0$$

$$H_1 : K(y) - 3 \neq 0.$$

- A estatística da razão *t*-Student da assimetria amostral é

$$t = \frac{\widehat{K}(y) - 3}{\sqrt{24/n}} \approx \mathcal{N}(0, 1).$$

- Rejeitamos H_0 ao nível α de significância se $|t| > z_{1-\alpha/2}$, em que $z_{1-\alpha/2}$ o percentil $100(1 - \alpha/2)$ da distribuição normal padrão.
- Temos também o teste de **Jarque e Bera** para normalidade com

$$JB = \frac{\widehat{A}^2(y)}{6/n} + \frac{(\widehat{K}(y) - 3)^2}{24/n} \approx \chi_2^2 \quad \text{para } n \text{ "grande"}.$$

- Rejeitamos H_0 (normalidade) se $JB > \chi_{1-\alpha}^*$, em que $\chi_{1-\alpha}^*$ o percentil $100(1 - \alpha)$ da distribuição χ_2^2 .

Função de autocorrelação (pausa em regressão)

- Considere uma sequência de observações (y_1, y_2, \dots, y_n) equiespaçadas na escala do tempo.
- Estamos interessados na correlação linear entre y_t e y_{t-h} para algum inteiro h .
- O coeficiente de correlação entre y_t e y_{t-h} é chamado de autocorrelação de defasagem (*lag*) h de y_t .

$$\rho(h) = \frac{\text{Cov}(y_t, y_{t-h})}{\sqrt{\text{Var}(y_t)\text{Var}(y_{t-h})}} = \frac{\text{Cov}(y_t, y_{t-h})}{\text{Var}(y_t)} = \frac{\gamma(h)}{\gamma(0)},$$

pois $\text{Var}(y_t) = \text{Var}(y_{t-h})$ (**hipótese de estacionariedade fraca**).

- Temos que $\rho(0) = 1$, $\rho(h) = \rho(-h)$ e $-1 \leq \rho(h) \leq 1$ para todo h .
- Uma sequência de observações (**fracamente estacionária**) y_t não é correlacionada serialmente se, e somente se, $\rho(h) = 0$ para todo $h > 0$.

Estimação da função de autocorrelação

- Para uma dada amostra $\{y_t\}_{t=1}^n$, a autocorrelação de defasagem 1 de y_t é

$$\hat{\rho}(1) = \frac{\sum_{t=2}^n (y_t - \bar{y})(y_{t-1} - \bar{y})}{\sum_{t=1}^n (y_t - \bar{y})^2}.$$

- Sob algumas condições gerais, $\hat{\rho}(1)$ é um estimador consistente de $\rho(1)$.
- Se $\{y_t\}$ for uma sequência independente e identicamente distribuída (i.i.d.) e $E(y_t^2) < \infty$, então $\hat{\rho}(1)$ é assintoticamente normal com média 0 e variância $1/n$.
- Para n suficientemente grande, temos

$$\hat{\rho}(1)\sqrt{n} \approx \mathcal{N}(0, 1).$$

- Podemos testar $H_0 : \rho(1) = 0$ contra $H_1 : \rho(1) \neq 0$.

Estimação da função de autocorrelação: defasagem $h \geq 0$

- A autocorrelação de defasagem h de y_t é definida por

$$\hat{\rho}(h) = \frac{\sum_{t=h+1}^n (y_t - \bar{y})(y_{t-h} - \bar{y})}{\sum_{t=1}^n (y_t - \bar{y})^2}, \quad \text{para } 0 \leq h \leq n-1.$$

- Se $\{y_t\}$ for uma sequência i.i.d. com $E(y_t^2) < \infty$, então $\hat{\rho}(h)$ é assintoticamente normal com média 0 e variância $1/n$ para todo inteiro positivo e fixo h .
- Para n suficientemente grande, temos

$$\hat{\rho}(h)\sqrt{n} \approx \mathcal{N}(0, 1).$$

- Podemos testar $H_0 : \rho(h) = 0$ contra $H_1 : \rho(h) \neq 0$ para h fixo.

Predição na média (de volta à regressão)

Dado uma regressão linear simples ajustada como

$$\hat{y} = b_0 + b_1 x = \bar{y} + b_1(x - \bar{x}),$$

para prever pontualmente o valor $E(y_0|\mathbf{x}, x_0) = \beta_0 + \beta_1 x_0$ no ponto x_0 , utilizamos

$$\hat{y}_0 = b_0 + b_1 x_0 = \bar{y} + b_1(x_0 - \bar{x}).$$

Além disso,

$$\begin{aligned}\text{Var}(\hat{y}_0|\mathbf{x}, x_0) &= \text{Var}(\bar{y}|\mathbf{x}) + (x_0 - \bar{x})^2 \text{Var}(b_1|\mathbf{x}) \\ &= \frac{\sigma_\varepsilon^2}{n} + \frac{(x_0 - \bar{x})^2 \sigma_\varepsilon^2}{S_{xx}} = \sigma_\varepsilon^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]; \quad e\end{aligned}$$

$$\text{DP}(\hat{y}_0|\mathbf{x}, x_0) = \sigma_\varepsilon \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]^{1/2}.$$

Consequentemente, temos que

$$\text{ep}(\hat{y}_0|\mathbf{x}, x_0) = \widehat{\text{DP}}(\hat{y}_0|\mathbf{x}, x_0) = s \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]^{1/2}.$$

Temos também que

$$\begin{aligned} E(\hat{y}_0 | \mathbf{x}, x_0) &= E(\bar{y} | \mathbf{x}) + (x_0 - \bar{x})E(b_1 | \mathbf{x}) \\ &= \beta_0 + \beta_1 \bar{x} + (x_0 - \bar{x})\beta_1 = \beta_0 + \beta_1 x_0, \end{aligned}$$

e, sob **HP.3**, $\hat{y}_0 = \bar{y} + b_1(x_0 - \bar{x})$ é uma combinação linear de normais. Logo,

$$(\hat{y}_0 | \mathbf{x}, x_0) \sim \mathcal{N} \left(\beta_0 + \beta_1 x_0, \sigma_\varepsilon^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right] \right)$$

é possível mostrar que s^2 e \hat{y}_0 são independentes. Daí, temos que

$$IC_{100(1-\alpha)\%} (E(y_0 | \mathbf{x}, x_0)) = (\hat{y}_0 - t_{n-2; 1-\alpha/2} \times \text{ep}(\hat{y}_0); \hat{y}_0 + t_{n-2; 1-\alpha/2} \times \text{ep}(\hat{y}_0)).$$

Ver arquivo exemplo_01.r

Predição para observações

Dado uma regressão linear simples $y = \beta_0 + \beta_1 x + \varepsilon$ e ajustada como

$$\hat{y} = b_0 + b_1 x = \bar{y} + b_1(x - \bar{x}),$$

para prever pontualmente o valor $y_0 = \beta_0 + \beta_1 x_0 + \varepsilon_0$ no ponto x_0 , utilizamos

$$\hat{y}_0 = b_0 + b_1 x_0 = \bar{y} + b_1(x_0 - \bar{x}).$$

Assim, o erro de previsão é dado por $\hat{\varepsilon}_0 = y_0 - \hat{y}_0$ e consequentemente

$$\begin{aligned}\text{Var}(\hat{\varepsilon}_0 | \mathbf{x}, x_0) &= \text{Var}(y_0 | \mathbf{x}, x_0) + \text{Var}(\hat{y}_0 | \mathbf{x}, x_0) \\ &= \sigma_\varepsilon^2 \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]; \text{ e} \\ \text{DP}(\hat{\varepsilon}_0 | \mathbf{x}, x_0) &= \sigma_\varepsilon \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]^{1/2}.\end{aligned}$$

Consequentemente, o erro padrão da previsão é dado por

$$\text{ep}(\hat{\varepsilon}_0 | \mathbf{x}, x_0) = \widehat{\text{DP}}(\hat{\varepsilon}_0 | \mathbf{x}, x_0) = s \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]^{1/2}.$$

Sob **HP.3**, \hat{y}_0 é combinação linear de normais. Logo,

$$(\hat{\varepsilon}_0 | \mathbf{x}, x_0) \sim \mathcal{N} \left(0, \sigma_\varepsilon^2 \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right] \right)$$

Daí, temos que

$$\text{IC}_{100(1-\alpha)\%}(y_0 | \mathbf{x}, x_0) = (\hat{y}_0 - t_{n-2; 1-\alpha/2} \times \text{ep}(\hat{\varepsilon}_0 | \mathbf{x}, x_0), \hat{y}_0 + t_{n-2; 1-\alpha/2} \times \text{ep}(\hat{\varepsilon}_0 | \mathbf{x}, x_0)).$$

Ver arquivo exemplo_01.r

Elasticidade

Em economia, o termo elasticidade se refere a como uma variável econômica muda em função de outra variável econômica.

A elasticidade é medida em termos percentuais ao invés de absoluto.

Isto significa que medimos uma mudança na variável como uma porcentagem da quantidade original da variável.

A mudança em porcentagem da variável x é definida por

$$\text{Mudança percentual em } x = \frac{\text{Mudança em } x}{\text{Valor original de } x} = \frac{\Delta x}{x} \quad (\text{de } x \text{ para } x + \Delta x).$$

A elasticidade de y com respeito a x é dada por

$$\epsilon = \frac{\text{Mudança percentual em } y}{\text{Mudança percentual em } x} = \frac{\Delta y/y}{\Delta x/x}.$$

Em termos infinitesimais, temos que

$$\epsilon = \frac{\partial y/y}{\partial x/x} = \frac{x}{y} \frac{\partial y}{\partial x}.$$

Modelo de regressão linear múltipla

Notação

$$\begin{aligned}y &= f(x_1, x_2, \dots, x_p) + \varepsilon \\ &= x_1\beta_1 + x_2\beta_2 + \dots + x_p\beta_p + \varepsilon,\end{aligned}$$

em que

- y é a variável dependente;
- $(x_1, x_2, \dots, x_p)^\top$ é o vetor de variáveis regressoras (explicativas);
- $(\beta_1, \beta_2, \dots, \beta_p)^\top$ é o vetor de coeficientes; e
- ε é um distúrbio aleatório.

Para o modelo de regressão acima, a **elasticidade de y com respeito à k -ésima variável x_k** é dada por

$$\epsilon = \frac{x_k}{y} \frac{\partial y}{\partial x_k} = \frac{x_k}{y} \beta_k.$$

Hipóteses do modelo de regressão linear

HP.1: Linearidade: $y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + \dots + x_{ip}\beta_p + \varepsilon_i$;

HP.2: Posto completo: não existe nenhuma relação exata entre as variáveis regressoras do modelo;

HP.3: Exogeneidade das variáveis regressoras:

$$E(\varepsilon_i | x_{i1}, x_{i2}, \dots, x_{ip}) = 0 \Rightarrow \text{Cov}(\mathbf{x}_i, \varepsilon_i) = \mathbf{0}$$

com $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})^\top$.

HP.4: Homocedasticidade e autocorrelação nula: cada distúrbio, ε_i , tem a mesma variância, σ^2 , e é não correlacionado com qualquer outro distúrbio ε_j ;

HP.5: Geração dos dados: os dados em $(x_{i1}, x_{i2}, \dots, x_{ip})$ podem ser qualquer mistura de constantes e variáveis aleatórias; e

HP.6: Distribuição normal: os distúrbios são normalmente distribuídos.

Notação

Temos que

$$\begin{aligned}y_i &= x_{i1}\beta_1 + x_{i2}\beta_2 + \cdots + x_{ip}\beta_p + \varepsilon_i \\ &= \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, \quad \text{para } i = 1, 2, \dots, n;\end{aligned}$$

- n é o número de observações;
- $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})^\top$ é o vetor da i -ésima observação das variáveis regressoras;
- $\mathbf{x}_k = (x_{k1}, x_{k2}, \dots, x_{kn})^\top$ é a amostra da k -ésima variável regressora;
- $\mathbf{y} = (y_1, y_2, \dots, y_n)^\top$ é a amostra da variável dependente;
- $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)^\top$ é o vetor dos distúrbios;
- $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_p)^\top$ é o vetor de coeficientes; e
- $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)$ é a matriz de variáveis regressoras.

Assim, o modelo de regressão linear múltipla pode ser reescrito por

$$\begin{aligned}\mathbf{y} &= \mathbf{x}_1\beta_1 + \mathbf{x}_2\beta_2 + \cdots + \mathbf{x}_p\beta_p + \boldsymbol{\varepsilon} \quad \text{ou} \\ \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.\end{aligned}$$

O modelo log-linear

$$\ln y = \beta_1 + \beta_2 \ln x_2 + \cdots + \beta_p \ln x_p + \varepsilon$$

(Elasticidade constante)

A **elasticidade de y com respeito a mudanças em x_k** é dado por

$$\epsilon = \frac{\partial y/y}{\partial x/x} = \frac{\partial y}{\partial x} \times \frac{\partial \ln y / \partial y}{\partial \ln x / \partial x} = \frac{\partial \ln y}{\partial \ln x_k} = \beta_k.$$

Exemplo - mercado de gasolina dos EUA - 1953-2004

Considere o seguinte modelo para o consumo de gasolina per capita:

$$\ln(G/Pop) = \beta_1 + \beta_2 \ln(Renda/Pop) + \beta_3 \ln \text{Preço} + \beta_4 \ln(PCN) + \beta_5 \ln(PCU) + \varepsilon,$$

sendo PCN o preço dos carros novos e PCU o preço dos carros usados.

- Elasticidade renda e preço para a gasolina; e
- Elasticidade para demanda com respeito a preço de carros novos e usados.

Distúrbios do modelo de regressão

Temos:

- $E(\varepsilon_i | \mathbf{x}_i) = 0$ por hipótese, ou seja,

$$E(\varepsilon | \mathbf{X}) = \begin{bmatrix} E(\varepsilon_1 | \mathbf{x}_1) \\ E(\varepsilon_2 | \mathbf{x}_2) \\ \vdots \\ E(\varepsilon_n | \mathbf{x}_n) \end{bmatrix} = \mathbf{0};$$

- $E(\varepsilon_i | \varepsilon_1, \dots, \varepsilon_{i-1}, \varepsilon_{i+1}, \dots, \varepsilon_n) = 0$;
- Por hipótese, todos os distúrbios são simplesmente uma amostra aleatória de alguma população;
- $E(\varepsilon_i) = E_{\mathbf{x}_i} (E(\varepsilon_i | \mathbf{x}_i)) = E_{\mathbf{x}_i} (0) = 0$;
- Para cada ε_i , $\text{Cov}(\varepsilon_i, \mathbf{x}_i) = \text{Cov}(E(\varepsilon_i | \mathbf{x}_i), \mathbf{x}_i) = \text{Cov}(0, \mathbf{x}_i) = \mathbf{0}$;

Note que $E(\varepsilon_i) = 0$ **não implica** que $E(\varepsilon_i | \mathbf{x}_i) = 0$. Isto é, $E(\varepsilon_i) = 0$ **não implica** que $\text{Cov}(\varepsilon_i, \mathbf{x}_i) = \mathbf{0}$.

- $\mathbf{y} = \mathbf{X}\beta + \varepsilon \Rightarrow E(\mathbf{y} | \mathbf{X}) = E(\mathbf{X}\beta | \mathbf{X}) + E(\varepsilon | \mathbf{X}) = \mathbf{X}\beta + \mathbf{0} = \mathbf{X}\beta$.

Distúrbios esféricos

Temos:

- $\text{Var}(\varepsilon_i | \mathbf{x}_i) = \sigma^2$, para todo $i = 1, 2, \dots, n$;
- $\text{Cov}(\varepsilon_i, \varepsilon_j | \mathbf{x}_i, \mathbf{x}_j) = 0$, para todo $i \neq j$ (correlação nula), ou seja,

$$\begin{aligned} E(\varepsilon\varepsilon^\top | \mathbf{X}) &= \begin{bmatrix} E(\varepsilon_1\varepsilon_1 | \mathbf{x}_1, \mathbf{x}_1) & E(\varepsilon_1\varepsilon_2 | \mathbf{x}_1, \mathbf{x}_2) & \cdots & E(\varepsilon_1\varepsilon_n | \mathbf{x}_1, \mathbf{x}_n) \\ E(\varepsilon_2\varepsilon_1 | \mathbf{x}_2, \mathbf{x}_1) & E(\varepsilon_2\varepsilon_2 | \mathbf{x}_2, \mathbf{x}_2) & \cdots & E(\varepsilon_2\varepsilon_n | \mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ E(\varepsilon_n\varepsilon_1 | \mathbf{x}_n, \mathbf{x}_1) & E(\varepsilon_n\varepsilon_2 | \mathbf{x}_n, \mathbf{x}_2) & \cdots & E(\varepsilon_n\varepsilon_n | \mathbf{x}_n, \mathbf{x}_n) \end{bmatrix} \\ &= \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{bmatrix} = \sigma^2 \mathbf{I}; \end{aligned}$$

- $\text{Var}(\varepsilon) = E_{\mathbf{X}}(\text{Var}(\varepsilon | \mathbf{X})) + \text{Var}_{\mathbf{X}}(E(\varepsilon | \mathbf{X})) = \sigma^2 \mathbf{I}$;
- Por hipótese, $(\varepsilon | \mathbf{X}) \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$.

Mínimos quadrados

Temos que

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i.$$

- $\boldsymbol{\beta}$ e ε_i são quantidades da população;
- A regressão populacional é $E(y_i|\mathbf{x}_i) = \mathbf{x}_i^\top \boldsymbol{\beta}$;
- \mathbf{b} e e_i são estimativas amostrais;
- A estimativa de $E(y_i|\mathbf{x}_i)$ é $\hat{y}_i = \mathbf{x}_i^\top \mathbf{b}$;
- Seja $\varepsilon_i = y_i - \mathbf{x}_i^\top \boldsymbol{\beta}$ o distúrbio da i -ésima observação;
- Para qualquer valor de \mathbf{b} , estimamos ε_i com o resíduo

$$e_i = y_i - \mathbf{x}_i^\top \mathbf{b};$$

- Das definições, $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i = \mathbf{x}_i^\top \mathbf{b} + e_i$;
- Podemos estimar $\boldsymbol{\beta}$ com a amostra de (y_i, \mathbf{x}_i) , $i = 1, 2, \dots, n$;
- Na estimação de $\boldsymbol{\beta}$ por \mathbf{b} , queremos que os pontos observados fiquem perto da reta ajustada;
- Perto neste caso é definido por algum critério de ajuste;
- **Critério de mínimos quadrados.**

Temos

$$S(\mathbf{b}) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{b})^2,$$

com \mathbf{b} denotando a escolha do vetor de coeficientes.

Queremos minimizar $S(\mathbf{b})$, ou seja,

$$\min_{\mathbf{b}} S(\mathbf{b}) = \min_{\mathbf{b}} \mathbf{e}^\top \mathbf{e} = \min_{\mathbf{b}} (\mathbf{y} - \mathbf{X}\mathbf{b})^\top (\mathbf{y} - \mathbf{X}\mathbf{b}).$$

Mas

$$\begin{aligned} S(\mathbf{b}) &= \mathbf{y}^\top \mathbf{y} - \mathbf{b}^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\mathbf{b} + \mathbf{b}^\top \mathbf{X}^\top \mathbf{X}\mathbf{b} \\ &= \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\mathbf{b} + \mathbf{b}^\top \mathbf{X}^\top \mathbf{X}\mathbf{b}. \end{aligned}$$

A condição necessária para um mínimo é

$$\frac{\partial S(\mathbf{b})}{\partial \mathbf{b}} = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\mathbf{b} = \mathbf{0}.$$

Seja \mathbf{b} a solução. Então, \mathbf{b} satisfaz as equações normais de mínimos quadrados:

$$\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{y}.$$

Se a inversa de $\mathbf{X}^T \mathbf{X}$ existir, que segue pela hipótese de posto completo de \mathbf{X} , então a solução é

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Para esta solução minimizar a soma de quadrados,

$$\frac{\partial^2 S(\mathbf{b})}{\partial \mathbf{b} \partial \mathbf{b}^T} = 2\mathbf{X}^T \mathbf{X}$$

deve ser uma matriz positiva definida.

Seja $q = \mathbf{c}^T \mathbf{X}^T \mathbf{X} \mathbf{c}$ para algum vetor arbitrário \mathbf{c} diferente de zero. Então,

$$q = \mathbf{v}^T \mathbf{v} = \sum_{i=1}^n v_i^2 \quad \text{sendo} \quad \mathbf{v} = \mathbf{X} \mathbf{c}.$$

A menos que cada elemento de \mathbf{v} seja zero, q é positivo.

Exemplo

$$y_i = \beta_1 + \beta_2 T_i + \beta_3 G_i + \beta_4 R_i + \beta_5 P_i + \varepsilon_i$$

Tabela 6: Observações anuais de 1968 a 1982.

Investimento Real (Y)	Constante (1)	Tendência (T)	PIB Real (G)	Taxa de Juros (R)	Taxa de Inflação (P)
0,161	1	1	1,058	5,16	4,40
0,172	1	2	1,088	5,87	5,15
0,158	1	3	1,086	5,95	5,37
0,173	1	4	1,122	4,88	4,99
0,195	1	5	1,186	4,50	4,16
0,217	1	6	1,254	6,44	5,75
0,199	1	7	1,246	7,83	8,82
0,163	1	8	1,232	6,25	9,31
0,195	1	9	1,298	5,50	5,21
0,231	1	10	1,370	5,46	5,83
0,257	1	11	1,439	7,46	7,40
0,259	1	12	1,479	10,28	8,64
0,225	1	13	1,474	11,77	9,31
0,241	1	14	1,503	13,42	9,44
0,204	1	15	1,475	11,02	5,99

Ver arquivo exemplo_07.r

Teorema (regressão ortogonal)

Se as variáveis em um modelo de regressão são não correlacionadas, isto é, são ortogonais, então os coeficientes angulares da regressão são os mesmos que os coeficientes na regressão individual simples.

Aspectos algébricos da solução de mínimos quadrados

- As equações normais são

$$\mathbf{X}^T \mathbf{X} \mathbf{b} - \mathbf{X}^T \mathbf{y} = -\mathbf{X}^T (\mathbf{y} - \mathbf{X} \mathbf{b}) = -\mathbf{X}^T \mathbf{e} = \mathbf{0}.$$

- Para cada coluna \mathbf{x}_k de \mathbf{X} , temos $\mathbf{x}_k^T \mathbf{e} = 0$.
- Se a primeira coluna de \mathbf{X} for de 1's, então
 1. os resíduos de mínimos quadrados somam zero,

$$\mathbf{x}_1^T \mathbf{e} = \mathbf{1}^T \mathbf{e} = \sum_{i=1}^n e_i = 0;$$

2. a regressão no hiperplano passa pelo ponto das médias dos dados; a primeira equação normal implica que $\bar{y} = \bar{\mathbf{x}}^T \mathbf{b}$;
3. a média dos valores ajustados da regressão é igual a média dos valores atuais.

Isto só vale se a regressão tiver uma constante.

Projeção

O valor dos resíduos de mínimos quadrados é dado por

$$\begin{aligned}\mathbf{e} &= \mathbf{y} - \mathbf{X}\mathbf{b} \\ &= \mathbf{y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{y} - \mathbf{P}\mathbf{y} \\ &= \mathbf{I}\mathbf{y} - \mathbf{P}\mathbf{y} \\ &= (\mathbf{I} - \mathbf{P})\mathbf{y} \\ &= \mathbf{M}\mathbf{y},\end{aligned}$$

sendo $\mathbf{P} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ e $\mathbf{M} = \mathbf{I} - \mathbf{P}$.

Temos que \mathbf{P} é a matriz de projeção (leva \mathbf{y} em $\hat{\mathbf{y}} = \mathbf{X}\mathbf{b}$):

$$\mathbf{P} = \mathbf{P}^\top, \quad \text{e} \quad \mathbf{P} = \mathbf{P}^2.$$

E \mathbf{M} é uma matriz importante, pois leva \mathbf{y} em \mathbf{e} :

$$\mathbf{M} = \mathbf{M}^\top, \quad \mathbf{M} = \mathbf{M}^2, \quad \text{e} \quad \mathbf{M}\mathbf{X} = \mathbf{0}.$$

Toda matriz simétrica e idempotente é positiva definida, em particular \mathbf{P} e \mathbf{M} .

A matriz de projeção

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e} \Leftrightarrow \mathbf{y} = \hat{\mathbf{y}} + \mathbf{e} \Rightarrow$$
$$\hat{\mathbf{y}} = \mathbf{y} - \mathbf{e} = (\mathbf{I} - \mathbf{M})\mathbf{y} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{P}\mathbf{y}.$$

$$\mathbf{P}\mathbf{M} = \mathbf{M}\mathbf{P} = \mathbf{0}, \quad \text{e} \quad \mathbf{P}\mathbf{X} = \mathbf{X}.$$

$$\mathbf{y} = \mathbf{P}\mathbf{y} + \mathbf{M}\mathbf{y} = \text{projeção} + \text{resíduos}.$$

Podemos reescrever a soma de quadrados da seguinte forma:

$$\begin{aligned}\mathbf{y}^T\mathbf{y} &= \mathbf{y}^T\mathbf{P}^T\mathbf{P}\mathbf{y} + \mathbf{y}^T\mathbf{M}^T\mathbf{M}\mathbf{y} \\ &= \hat{\mathbf{y}}^T\hat{\mathbf{y}} + \mathbf{e}^T\mathbf{e}.\end{aligned}$$

Outras relações úteis:

$$\begin{aligned}\mathbf{e}^T\mathbf{e} &= \mathbf{y}^T\mathbf{M}^T\mathbf{M}\mathbf{y} = \mathbf{y}^T\mathbf{M}\mathbf{y} = \mathbf{y}^T\mathbf{e} = \mathbf{e}^T\mathbf{y} \\ \mathbf{e}^T\mathbf{e} &= \mathbf{y}^T\mathbf{y} - \mathbf{b}^T\mathbf{X}^T\mathbf{X}\mathbf{b} = \mathbf{y}^T\mathbf{y} - \mathbf{b}^T\mathbf{X}^T\mathbf{y} = \mathbf{y}^T\mathbf{y} - \mathbf{y}^T\mathbf{X}\mathbf{b}.\end{aligned}$$

Regressão particionada e regressão parcial

Suponha que a regressão envolva dois conjuntos de variáveis \mathbf{X}_1 e \mathbf{X}_2 . Logo,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}.$$

Qual é a solução algébrica para \mathbf{b}_2 ?

$$\begin{bmatrix} \mathbf{X}_1^\top \mathbf{X}_1 & \mathbf{X}_1^\top \mathbf{X}_2 \\ \mathbf{X}_2^\top \mathbf{X}_1 & \mathbf{X}_2^\top \mathbf{X}_2 \end{bmatrix} \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1^\top \mathbf{y} \\ \mathbf{X}_2^\top \mathbf{y} \end{bmatrix}$$

Resolvendo para \mathbf{b}_1 (primeira equação):

$$\begin{aligned} \mathbf{X}_1^\top \mathbf{X}_1 \mathbf{b}_1 + \mathbf{X}_1^\top \mathbf{X}_2 \mathbf{b}_2 &= \mathbf{X}_1^\top \mathbf{y} \\ (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_1 \mathbf{b}_1 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \mathbf{b}_2 &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y} \end{aligned}$$

$$\begin{aligned} \mathbf{b}_1 &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y} - (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \mathbf{b}_2 \\ &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top (\mathbf{y} - \mathbf{X}_2 \mathbf{b}_2). \end{aligned}$$

\mathbf{b}_1 é o conjunto de coeficientes na regressão de \mathbf{y} sobre \mathbf{X}_1 , menos um vetor de correção. Suponha que $\mathbf{X}_1^\top \mathbf{X}_2 = \mathbf{0}$, então $\mathbf{b}_1 = (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y}$.

Na solução de \mathbf{b}_2 temos

$$(\mathbf{X}_2^\top \mathbf{X}_1) \mathbf{b}_1 + (\mathbf{X}_2^\top \mathbf{X}_2) \mathbf{b}_2 = \mathbf{X}_2^\top \mathbf{y}$$

$$(\mathbf{X}_2^\top \mathbf{X}_1)(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top (\mathbf{y} - \mathbf{X}_2 \mathbf{b}_2) + (\mathbf{X}_2^\top \mathbf{X}_2) \mathbf{b}_2 = \mathbf{X}_2^\top \mathbf{y}$$

$$\mathbf{X}_2^\top \mathbf{P}_1 \mathbf{y} - \mathbf{X}_2^\top \mathbf{P}_1 \mathbf{X}_2 \mathbf{b}_2 + (\mathbf{X}_2^\top \mathbf{X}_2) \mathbf{b}_2 = \mathbf{X}_2^\top \mathbf{y}$$

$$\mathbf{X}_2^\top (\mathbf{I} - \mathbf{P}_1) \mathbf{X}_2 \mathbf{b}_2 = \mathbf{X}_2^\top (\mathbf{I} - \mathbf{P}_1) \mathbf{y}$$

$$\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2 \mathbf{b}_2 = \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{y}$$

$$\mathbf{b}_2 = (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{y} = (\mathbf{X}_2^{*\top} \mathbf{X}_2^*)^{-1} \mathbf{X}_2^{*\top} \mathbf{y}^*,$$

sendo $\mathbf{X}_2^* = \mathbf{M}_1 \mathbf{X}_2$ e $\mathbf{y}^* = \mathbf{M}_1 \mathbf{y}$.

\mathbf{M}_1 é a matriz construtora de resíduos definida para a regressão sobre as colunas de \mathbf{X}_1 .

Logo, $\mathbf{M}_1 \mathbf{X}_2$ é uma matriz de resíduos.

Cada coluna de $\mathbf{M}_1 \mathbf{X}_2$ é um vetor de resíduos da regressão correspondente da coluna \mathbf{X}_2 sobre as variáveis em \mathbf{X}_1 .

Teorema (Frisch-Waugh, 1933; e Lovell, 1963)

Na regressão linear de mínimos quadrados do vetor \mathbf{y} sobre dois conjuntos de variáveis \mathbf{X}_1 e \mathbf{X}_2 , o subvetor \mathbf{b}_2 é o conjunto de coeficientes obtidos quando os resíduos da regressão de \mathbf{y} sobre \mathbf{X}_1 sozinho são regredidos sobre o conjunto de resíduos obtidos quando cada coluna de \mathbf{X}_2 é regredida sobre \mathbf{X}_1 .

Considere uma regressão de \mathbf{y} sobre um conjunto de variáveis \mathbf{X} e uma variável adicional \mathbf{z} , e seus respectivos coeficientes por \mathbf{b} e c .

Corolário

O coeficiente sobre \mathbf{z} em uma regressão múltipla de \mathbf{y} sobre $\mathbf{W} = [\mathbf{X}, \mathbf{z}]$ é calculada como

$$c = (\mathbf{z}^\top \mathbf{Mz})^{-1} \mathbf{z}^\top \mathbf{My} = (\mathbf{z}^{*\top} \mathbf{z}^*)^{-1} \mathbf{z}^* \mathbf{y}^*,$$

sendo $\mathbf{z}^* = \mathbf{Mz}$ e $\mathbf{y}^* = \mathbf{My}$ os vetores de resíduos da regressão de mínimos quadrados de \mathbf{z} e \mathbf{y} sobre \mathbf{X} .

Lembrando que $\mathbf{M} = \mathbf{I} - \mathbf{P} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$.

Coeficiente de correlação parcial

Seja a regressão $y = \beta_1 + \beta_2 z + \beta_3 x + \varepsilon$. Além disso,

(a) y_i^* os resíduos da regressão de y sobre x , isto é, $y_i^* = y_i - a - bx_i$; e

(b) z_i^* os resíduos da regressão de z sobre x , isto é, $z_i^* = z_i - c - dx_i$,

para $i = 1, 2, \dots, n$.

Então, a correlação parcial r_{yz}^* é a correlação simples entre y^* e z^* , isto é,

$$r_{yz}^* = \frac{\mathbf{z}^{*\top} \mathbf{y}^*}{\sqrt{(\mathbf{z}^{*\top} \mathbf{z}^*)(\mathbf{y}^{*\top} \mathbf{y}^*)}}.$$

pois o modelo de regressão tem a constante β_1 . Daí, $\mathbf{1}^\top \mathbf{z}^* = \mathbf{1}^\top \mathbf{y}^* = 0$.

Notemos que $\mathbf{y}^* = \mathbf{M}\mathbf{y}$ e $\mathbf{z}^* = \mathbf{M}\mathbf{z}$ com $\mathbf{M} = \mathbf{I} - \mathbf{P} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$.

Soma de quadrados

Considere agora os resíduos \mathbf{u} da regressão: $\mathbf{y} = \mathbf{X}\mathbf{d} + \mathbf{z}\mathbf{c} + \mathbf{u}$.

A menos que $\mathbf{X}^\top \mathbf{z} = \mathbf{0}$, \mathbf{d} não será igual a $\mathbf{b} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.

A menos que $\mathbf{c} = 0$, \mathbf{u} não será igual a $\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{b}$.

Notemos que

$$\begin{aligned}c &= (\mathbf{z}^{*\top} \mathbf{z}^*)^{-1} \mathbf{z}^{*\top} \mathbf{y}^*, \quad \mathbf{e} \\ \mathbf{d} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{z}c) = \mathbf{b} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{z}c.\end{aligned}$$

Segue-se que

$$\mathbf{u} = \mathbf{y} - \mathbf{X}\mathbf{b} + \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{z}c - \mathbf{z}c = \mathbf{e} - \mathbf{M}\mathbf{z}c = \mathbf{e} - \mathbf{z}^*c.$$

Portanto,

$$\mathbf{u}^\top \mathbf{u} = \mathbf{e}^\top \mathbf{e} + c^2(\mathbf{z}^{*\top} \mathbf{z}^*) - 2c\mathbf{z}^{*\top} \mathbf{e}.$$

Mas

$$\mathbf{e} = \mathbf{M}\mathbf{y} = \mathbf{y}^* \Rightarrow \mathbf{z}^{*\top} \mathbf{e} = \mathbf{z}^{*\top} \mathbf{y}^* = c(\mathbf{z}^{*\top} \mathbf{z}^*),$$

pois pela equação normal $(\mathbf{z}^{*\top} \mathbf{z}^*)c = \mathbf{z}^{*\top} \mathbf{y}^*$.

Assim, $\mathbf{u}^\top \mathbf{u} = \mathbf{e}^\top \mathbf{e} - c^2(\mathbf{z}^{*\top} \mathbf{z}^*)$.

Teorema (mudança na soma de quadrados quando uma variável é adicionada ao modelo)

Se $\mathbf{e}^\top \mathbf{e}$ é a soma de quadrados dos resíduos quando \mathbf{y} é regredido sobre \mathbf{X} e $\mathbf{u}^\top \mathbf{u}$ é a soma de quadrados dos resíduos quando \mathbf{y} é regredido sobre \mathbf{X} e \mathbf{z} , então

$$\mathbf{u}^\top \mathbf{u} = \mathbf{e}^\top \mathbf{e} - c^2(\mathbf{z}^{*\top} \mathbf{z}^*) \leq \mathbf{e}^\top \mathbf{e},$$

sendo c o coeficiente de \mathbf{z} na regressão completa e $\mathbf{z}^* = [\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top] \mathbf{z}$ o vetor de resíduos quando \mathbf{z} é regredido sobre \mathbf{X} .

Temos que $\mathbf{e}^\top \mathbf{e} = \mathbf{y}^{*\top} \mathbf{y}^*$ e pela equação normal

$$c(\mathbf{z}^{*\top} \mathbf{z}^*) = \mathbf{z}^{*\top} \mathbf{y}^* \quad \Rightarrow \quad c^2(\mathbf{z}^{*\top} \mathbf{z}^*) = (\mathbf{z}^{*\top} \mathbf{y}^*)^2 / (\mathbf{z}^{*\top} \mathbf{z}^*).$$

Portanto,

$$\frac{\mathbf{u}^\top \mathbf{u}}{\mathbf{y}^{*\top} \mathbf{y}^*} = \frac{\mathbf{y}^{*\top} \mathbf{y}^* - (\mathbf{z}^{*\top} \mathbf{y}^*)^2 / \mathbf{z}^{*\top} \mathbf{z}^*}{\mathbf{y}^{*\top} \mathbf{y}^*} = 1 - r_{yz}^{*2}.$$

Exemplo (continuação)

$$y_i = \beta_1 + \beta_2 T_i + \beta_3 G_i + \beta_4 R_i + \beta_5 P_i + \varepsilon_i$$

Seguimos a análise dos dados da equação de investimento.

Calculamos a correlação (simples) entre a variável investimento e as variáveis de tendência temporal (T), produto interno bruto (G), taxa de juros (R) e taxa de inflação (P).

Calculamos também a correlação parcial entre investimento e as quatro regressoras dado as demais regressoras.

Tabela 7: Correlação entre investimento e as demais variáveis.

	Correlação Simples	Correlação Parcial
Tendência	0,7496	-0,9360
PIB	0,8632	0,9680
Taxa de juros	0,5871	-0,5167
Taxa de inflação	0,4777	-0,0221

Ver arquivo exemplo_07.r

Qualidade de ajuste e análise de variância

Podemos perguntar se a variação em \mathbf{x} é um bom preditor da variação em y .

A variação total é dada por:

$$SQTot = \sum_{i=1}^n (y_i - \bar{y})^2.$$

Já sabemos que $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e} = \hat{\mathbf{y}} + \mathbf{e}$. Logo, para uma única observação, temos que

$$\begin{aligned} y_i &= \hat{y}_i + e_i = \mathbf{x}_i^T \mathbf{b} + e_i \\ y_i - \bar{y} &= \hat{y}_i - \bar{y} + e_i = (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{b} + e_i. \end{aligned}$$

Algumas propriedades, notações e somas de quadrados

$$\mathbf{1}\bar{x} = \mathbf{1}\frac{1}{n}\mathbf{1}^T \mathbf{x} = \begin{bmatrix} \bar{x} \\ \bar{x} \\ \vdots \\ \bar{x} \end{bmatrix} = \frac{1}{n}\mathbf{1}\mathbf{1}^T \mathbf{x}$$

$$\begin{bmatrix} x_1 - \bar{x} \\ x_2 - \bar{x} \\ \vdots \\ x_n - \bar{x} \end{bmatrix} = [\mathbf{x} - \mathbf{1}\bar{x}] = [\mathbf{x} - \frac{1}{n}\mathbf{1}\mathbf{1}^T \mathbf{x}] = [I\mathbf{x} - \frac{1}{n}\mathbf{1}\mathbf{1}^T \mathbf{x}] = [I - \frac{1}{n}\mathbf{1}\mathbf{1}^T]\mathbf{x} = \tilde{\mathbf{M}}\mathbf{x}.$$

Definimos $\tilde{\mathbf{M}} = I - \frac{1}{n}\mathbf{1}\mathbf{1}^T$. Esta também é uma matriz idempotente.

Temos que $\tilde{\mathbf{M}}\mathbf{1} = \mathbf{0} \Leftrightarrow \mathbf{1}^T \tilde{\mathbf{M}} = \mathbf{0}^T$. Logo,

$$\sum_{i=1}^n (x_i - \bar{x}) = \mathbf{1}^T [\tilde{\mathbf{M}}\mathbf{x}] = \mathbf{0}^T \mathbf{x} = 0.$$

Além disso,

$$\sum_{i=1}^n (x_i - \bar{x})^2 = (\mathbf{x} - \mathbf{1}\bar{x})^T (\mathbf{x} - \mathbf{1}\bar{x}) = (\tilde{\mathbf{M}}\mathbf{x})^T (\tilde{\mathbf{M}}\mathbf{x}) = \mathbf{x}^T \tilde{\mathbf{M}}^T \tilde{\mathbf{M}}\mathbf{x} = \mathbf{x}^T \tilde{\mathbf{M}}\mathbf{x}$$

$$\begin{aligned}\tilde{\mathbf{M}}\mathbf{y} &= \tilde{\mathbf{M}}\mathbf{X}\mathbf{b} + \tilde{\mathbf{M}}\mathbf{e} \\ \tilde{\mathbf{M}}\mathbf{e} &= \mathbf{e} \quad \mathbf{e} \quad \mathbf{e}^\top \tilde{\mathbf{M}}\mathbf{X} = \mathbf{e}^\top \mathbf{X} = \mathbf{0}^\top.\end{aligned}$$

Finalmente,

$$\underbrace{\mathbf{y}^\top \tilde{\mathbf{M}}\mathbf{y}}_{SQ_{Tot}} = \underbrace{\mathbf{b}^\top \mathbf{X}^\top \tilde{\mathbf{M}}\mathbf{X}\mathbf{b}}_{SQ_{Reg}} + \underbrace{\mathbf{e}^\top \mathbf{e}}_{SQ_{Res}}.$$

- SQ_{Tot} é a soma de quadrados total;
- SQ_{Reg} é a soma de quadrados da regressão; e
- SQ_{Res} é a soma de quadrados dos resíduos.

Coeficiente de determinação

$$R^2 = \frac{SQ_{Reg}}{SQ_{Tot}} = \frac{\mathbf{b}^\top \mathbf{X}^\top \tilde{\mathbf{M}}\mathbf{X}\mathbf{b}}{\mathbf{y}^\top \tilde{\mathbf{M}}\mathbf{y}} = 1 - \frac{\mathbf{e}^\top \mathbf{e}}{\mathbf{y}^\top \tilde{\mathbf{M}}\mathbf{y}}.$$

Tabela 8: Análise de variância.

Soma de Quadrados	Fonte	Graus de Liberdade	Média Quadrática
Regressão	$\mathbf{b}^\top \mathbf{X}^\top \mathbf{y} - n\bar{y}^2$	$p - 1$	s^2
Resíduo	$\mathbf{e}^\top \mathbf{e}$	$n - p$	
Total	$\mathbf{y}^\top \mathbf{y} - n\bar{y}^2$	$n - 1$	$S_{yy}/(n - 1) = s_y^2$

$$\begin{aligned}
 \mathbf{b}^\top \mathbf{X}^\top \tilde{\mathbf{M}} \mathbf{X} \mathbf{b} &= \mathbf{b}^\top \mathbf{X}^\top \tilde{\mathbf{M}} [\mathbf{y} - \mathbf{e}] = \mathbf{b}^\top \mathbf{X}^\top \tilde{\mathbf{M}} \mathbf{y} - \mathbf{b}^\top \mathbf{X}^\top \tilde{\mathbf{M}} \mathbf{e} \\
 &= \mathbf{b}^\top \mathbf{X}^\top \left[\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right] \mathbf{y} \quad (\text{porque } \mathbf{X}^\top \tilde{\mathbf{M}} \mathbf{e} = \mathbf{0}) \\
 &= \mathbf{b}^\top \mathbf{X}^\top \mathbf{y} - \frac{1}{n} \mathbf{b}^\top \mathbf{X}^\top \mathbf{1} \mathbf{1}^\top \mathbf{y} = \mathbf{b}^\top \mathbf{X}^\top \mathbf{y} - \frac{1}{n} \hat{\mathbf{y}}^\top \mathbf{1} \mathbf{1}^\top \mathbf{y} \\
 &= \mathbf{b}^\top \mathbf{X}^\top \mathbf{y} - \frac{1}{n} (\mathbf{y} - \mathbf{e}) \mathbf{1} \mathbf{1}^\top \mathbf{y} \\
 &= \mathbf{b}^\top \mathbf{X}^\top \mathbf{y} - \frac{1}{n} \left[\sum_{i=1}^n y_i - \sum_{i=1}^n e_i \right] \sum_{i=1}^n y_i \\
 &= \mathbf{b}^\top \mathbf{X}^\top \mathbf{y} - n\bar{y}^2.
 \end{aligned}$$

Exemplo

$$y_i = \beta_1 + \beta_2 T_i + \beta_3 G_i + \beta_4 R_i + \beta_5 P_i + \varepsilon_i$$

Tabela 9: Análise de variância - equação de investimento.

Soma de Quadrados	Fonte	Graus de Liberdade	Média Quadrática	R^2
Regressão	0,0159025	4	0,00397563	0,9724
Resíduo	0,0004508	10	0,00004508	
Total	0,0163533	14	0,00116810	

Ver arquivo exemplo_07.r

Teorema (mudança no R^2 quando uma variável é adicionada ao modelo)

Seja R_{Xz}^2 o coeficiente de determinação na regressão de \mathbf{y} sobre \mathbf{X} e uma variável adicional \mathbf{z} (**\mathbf{u} os resíduos**), seja R_X^2 o mesmo para a regressão de \mathbf{y} sobre \mathbf{X} sozinha (**\mathbf{e} os resíduos**), e seja r_{yz}^{*2} o coeficiente de correlação parcial entre y e z . Então,

$$R_{Xz}^2 = R_X^2 + (1 - R_X^2)r_{yz}^{*2} \geq R_X^2.$$

Prova: Vimos que $\mathbf{z}^* = \mathbf{Mz}$, $\mathbf{y}^* = \mathbf{My} = \mathbf{e} \Rightarrow \mathbf{e}^\top \mathbf{e} = \mathbf{y}^{*\top} \mathbf{y}^*$, e

$$\mathbf{u}^\top \mathbf{u} = \mathbf{e}^\top \mathbf{e} - c^2(\mathbf{z}^{*\top} \mathbf{z}^*) = \mathbf{y}^{*\top} \mathbf{y}^* - \frac{(\mathbf{z}^{*\top} \mathbf{y}^*)^2}{\mathbf{z}^{*\top} \mathbf{z}^*} = \mathbf{e}^\top \mathbf{e}(1 - r_{yz}^{*2}).$$

Logo,

$$R_{Xz}^2 = 1 - \frac{\mathbf{u}^\top \mathbf{u}}{\mathbf{y}^\top \tilde{\mathbf{M}}\mathbf{y}} = 1 - \frac{\mathbf{e}^\top \mathbf{e}}{\mathbf{y}^\top \tilde{\mathbf{M}}\mathbf{y}} + \frac{\mathbf{e}^\top \mathbf{e}}{\mathbf{y}^\top \tilde{\mathbf{M}}\mathbf{y}} r_{yz}^{*2} = R_X^2 + (1 - R_X^2)r_{yz}^{*2} \geq R_X^2.$$

Exemplo - ajuste de uma função consumo

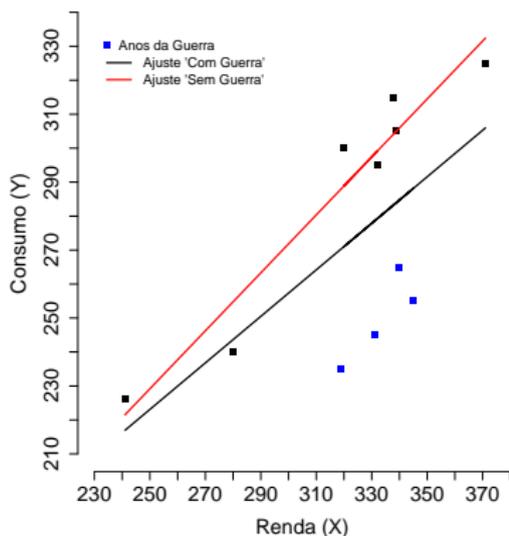


Figura 8: Dados de consumo e renda de 1940-1950 dos EUA. Fonte: *Economic Report of the President, U.S. Government Printing Office, Washington, D.C., 1983.*

Ver arquivo exemplo_08.r

Exemplo (continuação) - ajuste de uma função consumo

Tabela 10: Análise de variância: consumo e renda.

Soma de Quadrados	Fonte	Graus de Liberdade	Média Quadrática	R^2
Regressão	5768,2068	1	5768,2068	0,4571
Resíduo	6849,9751	9	761,1083	
Total	12618,1818	10	1261,8182	

Tabela 11: Análise de variância: consumo, renda e dicotômica para anos de guerra.

Soma de Quadrados	Fonte	Graus de Liberdade	Média Quadrática	R^2
Regressão x	5768,2068	1	5768,2068	0,9464
Regressão w	6173,5189	1	6173,5189	
Regressão	11941,7257	2	5970,8629	
Resíduo	676,4562	8	84,5570	
Total	12618,1818	10	1261,8,182	

O R^2 ajustado para os graus de liberdade

R^2 nunca decresce quando uma nova variável regressora é adicionado a equação de regressão. Portanto, é tentador adicionar diversas variáveis regressoras no modelo para que R^2 se aproxime de 1.

O R^2 **ajustado** para os graus de liberdade é dado por

$$\bar{R}^2 = 1 - \frac{\mathbf{e}^\top \mathbf{e} / (n - p)}{\mathbf{y}^\top \tilde{\mathbf{M}} \mathbf{y} / (n - 1)}.$$

Para efeitos computacionais, temos

$$\bar{R}^2 = 1 - \frac{n - 1}{n - p} \times (1 - R^2).$$

- O \bar{R}^2 **pode decrescer** quando uma nova covariável é adicionada ao modelo.
- Na verdade, \bar{R}^2 **pode até ser negativo**.

R^2 e o termo constante do modelo

A prova que $0 \leq R^2 \leq 1$ requer que \mathbf{X} contenha uma coluna de 1's.

Se isto não ocorrer, então $\tilde{\mathbf{M}}\mathbf{e} \neq \mathbf{e}$, $\mathbf{e}^\top \tilde{\mathbf{M}}\mathbf{X} \neq \mathbf{0}$ e o termo $2\mathbf{e}^\top \tilde{\mathbf{M}}\mathbf{X}\mathbf{b}$ em

$$\mathbf{y}^\top \tilde{\mathbf{M}}\mathbf{y} = (\tilde{\mathbf{M}}\mathbf{X}\mathbf{b} + \tilde{\mathbf{M}}\mathbf{e})^\top (\tilde{\mathbf{M}}\mathbf{X}\mathbf{b} + \tilde{\mathbf{M}}\mathbf{e})$$

não se cancela.

Daí,

$$R^2 = 1 - \frac{\mathbf{e}^\top \mathbf{e}}{\mathbf{y}^\top \tilde{\mathbf{M}}\mathbf{y}}$$

pode resultar em qualquer valor.

Portanto, devemos ter cuidado ao analisar o R^2 de um modelo de regressão que não contenha o termo constante.

Este é um dos motivos que, em geral, deixa-se o termo constante no modelo de regressão linear mesmo que testes de hipóteses indiquem o contrário.

Regressão transformada linearmente

Na regressão de \mathbf{y} sobre \mathbf{X} , suponha que as colunas de \mathbf{X} são transformadas linearmente.

Exemplo - precificação de arte

A teoria I da determinação de preços de leilão em pinturas de Monet diz que o preço é determinado pelas dimensões (largura W e altura H) da pintura,

$$\begin{aligned}\ln \text{preço} &= \beta_1 + \beta_2 \ln W + \beta_3 \ln H + \varepsilon \\ &= \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon.\end{aligned}$$

A teoria II diz que os compradores estão interessados na área da pintura e a razão do aspecto,

$$\begin{aligned}\ln \text{preço} &= \gamma_1 + \gamma_2 \ln(WH) + \gamma_3 \ln(W/H) + \xi \\ &= \gamma_1 z_1 + \gamma_2 z_2 + \gamma_3 z_3 + \xi.\end{aligned}$$

É evidente que $z_1 = x_1$, $z_2 = x_2 + x_3$ e $z_3 = x_2 - x_3$.

Em termos matriciais, temos $\mathbf{Z} = \mathbf{XP}$ sendo

$$\mathbf{P} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & -1 \end{bmatrix}.$$

Teorema

Na regressão linear de \mathbf{y} sobre $\mathbf{Z} = \mathbf{XP}$ sendo \mathbf{P} uma matriz não singular que transforma as colunas de \mathbf{X} , os coeficientes são iguais a $\mathbf{P}^{-1}\mathbf{b}$ sendo \mathbf{b} o vetor de coeficientes da regressão linear de \mathbf{y} sobre \mathbf{X} . Além disto, os R^2 das duas regressões são idênticos.

Prova: Os coeficientes são

$$\begin{aligned} \mathbf{d} &= (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y} = [(\mathbf{XP})^T (\mathbf{XP})]^{-1} (\mathbf{XP})^T \mathbf{y} = (\mathbf{P}^T \mathbf{X}^T \mathbf{XP})^{-1} \mathbf{P}^T \mathbf{X}^T \mathbf{y} \\ &= \mathbf{P}^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{P}^T \mathbf{P}^T \mathbf{X}^T \mathbf{y} = \mathbf{P}^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{P}^{-1} \mathbf{b}. \end{aligned}$$

Temos também que

$$\mathbf{u} = \mathbf{y} - \mathbf{Zd} = \mathbf{y} - \mathbf{ZP}^{-1}\mathbf{b} = \mathbf{y} - \mathbf{XPP}^{-1}\mathbf{b} = \mathbf{y} - \mathbf{Xb} = \mathbf{e}.$$

Logo, $\mathbf{u}^T \mathbf{u} = \mathbf{e}^T \mathbf{e} \Rightarrow R_{y|z}^2 = R_{y|x}^2$.

Hipóteses do modelo de regressão linear

- HP.1:** Linearidade: $y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + \dots + x_{ip}\beta_p + \varepsilon_i$;
- HP.2:** Posto completo: a matriz \mathbf{X} ($n \times p$) é de posto completo;
- HP.3:** Exogeneidade das variáveis independentes: $E(\varepsilon_i | \mathbf{x}_i) = 0$;
- HP.4:** Homocedasticidade e autocorrelação nula: cada distúrbio, ε_i , tem a mesma variância, σ^2 , e é não correlacionado com qualquer outro distúrbio ε_j ;
- HP.5:** Geração dos dados: os dados em $(x_{i1}, x_{i2}, \dots, x_{ip})$ podem ser qualquer mistura de constantes e variáveis aleatórias; e
- HP.6:** Distribuição normal: os distúrbios são normalmente distribuídos.

O estimador de mínimos quadrados

Por hipótese,

$$E(\varepsilon|\mathbf{x}) = 0 \Rightarrow \text{Cov}(\mathbf{x}, \varepsilon) = \text{Cov}(E(\varepsilon|\mathbf{x}), \mathbf{x}) = \mathbf{0} \quad \text{e} \quad E(\varepsilon) = E_{\mathbf{x}}(E(\varepsilon|\mathbf{x})) = 0.$$

Mas,

$$\begin{aligned} \text{Cov}(\mathbf{x}, \varepsilon) &= E_{\mathbf{x}, \varepsilon}(\mathbf{x}\varepsilon) = E_{\mathbf{x}, y}(\mathbf{x}(y - \mathbf{x}^T\boldsymbol{\beta})) = \mathbf{0} \Rightarrow \\ E_{\mathbf{x}, y}(\mathbf{x}y) &= E_{\mathbf{x}}(\mathbf{x}\mathbf{x}^T)\boldsymbol{\beta} \quad \text{(população)}. \end{aligned}$$

Tomemos a equação normal de mínimos quadrados $\mathbf{X}^T\mathbf{y} = \mathbf{X}^T\mathbf{X}\mathbf{b}$ que dividida por n , podemos reescrever como

$$\left(\frac{1}{n}\sum_{i=1}^n \mathbf{x}_i y_i\right) = \left(\frac{1}{n}\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T\right) \mathbf{b} \quad \text{(amostra)}.$$

Assim, o estimador de mínimos está parecido com sua respectiva quantidade populacional. Podemos aplicar a lei dos grandes números sob certas condições.

Preditor linear ótimo

Considere o problema de encontrar um preditor linear ótimo para y .

Critério: o erro quadrático médio (EQM) do preditor linear y .

Ignoraremos a **HP.6** e desconsideraremos **HP.1** que $E(y|\mathbf{x})$ é linear.

Queremos minimizar

$$\text{EQM} = E_{\mathbf{x},y}([y - \mathbf{x}^T \boldsymbol{\gamma}]^2).$$

Isto pode ser reescrito como

$$\text{EQM} = E_{\mathbf{x},y}([y - E(y|\mathbf{x})]^2) + E_{\mathbf{x},y}([E(y|\mathbf{x}) - \mathbf{x}^T \boldsymbol{\gamma}]^2).$$

A condição necessária é

$$\begin{aligned} \frac{\partial E_{\mathbf{x},y}([E(y|\mathbf{x}) - \mathbf{x}^T \boldsymbol{\gamma}]^2)}{\partial \boldsymbol{\gamma}} &= E_{\mathbf{x},y} \left(\frac{\partial [E(y|\mathbf{x}) - \mathbf{x}^T \boldsymbol{\gamma}]^2}{\partial \boldsymbol{\gamma}} \right) \\ &= -2E_{\mathbf{x},y}(\mathbf{x}[E(y|\mathbf{x}) - \mathbf{x}^T \boldsymbol{\gamma}]) = \mathbf{0} \Rightarrow \\ E_{\mathbf{x},y}(\mathbf{x}E(y|\mathbf{x})) &= E_{\mathbf{x},y}(\mathbf{x}\mathbf{x}^T)\boldsymbol{\gamma}. \end{aligned}$$

Mas

$$\begin{aligned} E_{\mathbf{x},y}(\mathbf{x}E(y|\mathbf{x})) &= \text{Cov}(\mathbf{x}, E(y|\mathbf{x})) + E_{\mathbf{x}}(\mathbf{x})E_{\mathbf{x}}(E(y|\mathbf{x})) \\ &= \text{Cov}(\mathbf{x}, y) + E_{\mathbf{x}}(\mathbf{x})E_y(y) = E_{\mathbf{x},y}(\mathbf{x}y). \end{aligned}$$

Portanto, a condição necessária para minimizar o preditor linear, EQM, é

$$E_{\mathbf{x},y}(\mathbf{x}y) = E_{\mathbf{x},y}(\mathbf{x}\mathbf{x}^T)\gamma.$$

Teorema

Se o processo de geração dos dados (y_i, \mathbf{x}_i) , $i = 1, 2, \dots, n$, é tal que a lei dos grandes números se aplica ao estimador

$$\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i y_i \right) = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{b}$$

das matrizes

$$E_{\mathbf{x},y}(\mathbf{x}y) = E_{\mathbf{x}}(\mathbf{x}\mathbf{x}^T)\beta,$$

então o mínimo do preditor linear do erro quadrático esperado de y é estimado pela regressão de mínimos quadrados.

Estimação não tendenciosa

O estimador de mínimos quadrados é não tendencioso para cada amostra.

Temos

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \boldsymbol{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} \quad \text{e}$$

$$E(\mathbf{b} | \mathbf{X}) = \boldsymbol{\beta} + E((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} | \mathbf{X}) = \boldsymbol{\beta} \quad (\text{pela HP.3}).$$

Portanto,

$$E(\mathbf{b}) = E_X(E(\mathbf{b} | \mathbf{X})) = E_X(\boldsymbol{\beta}) = \boldsymbol{\beta}.$$

Viés: omissão de variáveis relevantes

Suponha que o modelo correto seja

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$$

sendo p_1 e p_2 os números de colunas em \mathbf{X}_1 e \mathbf{X}_2 , respectivamente.

A regressão de \mathbf{y} sobre \mathbf{X}_1 (sem incluir \mathbf{X}_2) resulta no estimador

$$\mathbf{b}_1 = (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y} = \boldsymbol{\beta}_1 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \boldsymbol{\beta}_2 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \boldsymbol{\varepsilon}.$$

Logo,

$$E(\mathbf{b}_1 | \mathbf{X}) = \boldsymbol{\beta}_1 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \boldsymbol{\beta}_2 = \boldsymbol{\beta}_1 + \mathbf{P}_{1.2} \boldsymbol{\beta}_2,$$

sendo $\mathbf{P}_{1.2} = (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2$.

Portanto, \mathbf{b}_1 é um estimador tendencioso de $\boldsymbol{\beta}_1$, exceto para os casos que $\mathbf{X}_1^\top \mathbf{X}_2 = \mathbf{0}$ ou $\boldsymbol{\beta}_2 = \mathbf{0}$.

A estimação do modelo pode ser vista ao impor **incorretamente** a restrição $\boldsymbol{\beta}_2 = \mathbf{0}$. Isto introduz um viés na estimação de \mathbf{b}_1 .

Inclusão de variáveis irrelevantes

Suponha que o modelo de regressão correto seja

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \varepsilon,$$

mas que estimamos como a modelo de regressão (variáveis extras)

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \varepsilon.$$

Aqui, a inclusão de \mathbf{X}_2 no modelo de regressão pode ser vista como uma falha em impor que $\boldsymbol{\beta}_2 = \mathbf{0}$. Portanto, não a nada a se provar aqui.

Temos que

$$E(\mathbf{b}|\mathbf{X}) = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \mathbf{0} \end{bmatrix}.$$

É possível provar que ao aumentar o número de variáveis desnecessárias no modelo faz com que a precisão das estimativas diminua (a menos que $\mathbf{X}_1^\top \mathbf{X}_2 = \mathbf{0}$).

Variância amostral do estimador de mínimos quadrados

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\beta + \varepsilon) = \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon.$$

Logo,

$$\begin{aligned} \text{Var}(\mathbf{b}|\mathbf{X}) &= E([\mathbf{b} - \beta][\mathbf{b} - \beta]^T | \mathbf{X}) \\ &= E((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon \varepsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} | \mathbf{X}) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E(\varepsilon \varepsilon^T | \mathbf{X}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \quad (\text{pela HP.4}) \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}. \end{aligned}$$

Suponha uma regressão linear simples da forma $y = \beta_1 + \beta_2 x + \varepsilon$. Assim, temos

$$\text{Var}(b_2 | \mathbf{X}) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Agora, seja $\mathbf{b}_* = \mathbf{C}\mathbf{y}$ um outro estimador não tendencioso e linear de β , em que \mathbf{C} é uma matriz $p \times n$. Então,

$$\beta = E(\mathbf{b}_*|\mathbf{X}) = E(\mathbf{C}\mathbf{y}|\mathbf{X}) = E(\mathbf{C}\mathbf{X}\beta + \mathbf{C}\varepsilon|\mathbf{X}) \Rightarrow \mathbf{C}\mathbf{X} = \mathbf{I}.$$

(existem vários candidatos para \mathbf{C} que satisfazem a igualdade) e

$$\text{Var}(\mathbf{b}_*|\mathbf{X}) = \sigma^2 \mathbf{C}\mathbf{C}^\top.$$

Agora, seja $\mathbf{D} = \mathbf{C} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ tal que $\mathbf{D}\mathbf{y} = \mathbf{b}_* - \mathbf{b}$. Então, $\mathbf{b}_* = \mathbf{D}\mathbf{y} + \mathbf{b} = (\mathbf{D} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)\mathbf{y}$ e

$$\text{Var}(\mathbf{b}_*|\mathbf{X}) = \sigma^2 \left[(\mathbf{D} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)(\mathbf{D} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top \right].$$

Sabemos que $\mathbf{C}\mathbf{X} = \mathbf{I} = \mathbf{D}\mathbf{X} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \Rightarrow \mathbf{D}\mathbf{X} = \mathbf{0}$. Portanto,

$$\text{Var}(\mathbf{b}_*|\mathbf{X}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} + \sigma^2 \mathbf{D}\mathbf{D}^\top = \text{Var}(\mathbf{b}|\mathbf{X}) + \sigma^2 \mathbf{D}\mathbf{D}^\top.$$

e $\mathbf{D}\mathbf{D}^\top$ é uma forma quadrática ($\mathbf{a}^\top \mathbf{D}\mathbf{D}^\top \mathbf{a} \geq 0$).

Logo, $\text{Var}(\mathbf{b}|\mathbf{X}) \leq \text{Var}(\mathbf{b}_*|\mathbf{X})$ para qualquer $\mathbf{b}_* \neq \mathbf{b}$ e \mathbf{X} fixo.

Teorema de Gauss-Markov

No modelo de regressão linear clássico com matriz de variáveis regressoras \mathbf{X} , o estimador de mínimos quadrados \mathbf{b} é o estimador não tendencioso e linear de variância mínima de β . Para qualquer vetor de constantes \mathbf{v} , o estimador não tendencioso e linear de variância mínima de $\mathbf{v}^\top \beta$ no modelo de regressão clássico é $\mathbf{v}^\top \mathbf{b}$, sendo \mathbf{b} o estimador de mínimos quadrados.

O teorema também implica que o estimador de mínimos quadrados para cada coeficiente também tem variância mínima. Basta tomar um vetor $\mathbf{v} = (0, \dots, 0, 1, 0, \dots, 0)$ com o valor 1 no coeficiente a ser estimado.

A variância incondicional é dada por

$$\begin{aligned}\text{Var}(\mathbf{b}) &= E_{\mathbf{X}}(\text{Var}(\mathbf{b}|\mathbf{X})) + \text{Var}_{\mathbf{X}}(E(\mathbf{b}|\mathbf{X})) \\ &= \sigma^2 E_{\mathbf{X}}((\mathbf{X}^\top \mathbf{X})^{-1}),\end{aligned}$$

pois $E(\mathbf{b}|\mathbf{X}) = \beta \Rightarrow \text{Var}_{\mathbf{X}}(E(\mathbf{b}|\mathbf{X})) = \mathbf{0}$.

Estimação da variância

Sabemos que $\text{Var}(\mathbf{b}|\mathbf{X}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$. Precisamos estimar σ^2 .

Temos que $\text{Var}(\varepsilon_i) = \text{E}(\varepsilon_i^2) = \sigma^2$ e e_i é uma estimativa de ε_i . Então,

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n e_i^2$$

parece um estimador natural de σ^2 .

Os resíduos de mínimos quadrados são

$$\mathbf{e} = \mathbf{M}\mathbf{y} = \mathbf{M}[\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}] = \mathbf{M}\boldsymbol{\varepsilon}, \quad \text{pois } \mathbf{M}\mathbf{X} = \mathbf{0}.$$

Este estimador de σ^2 é baseado em $\mathbf{e}^\top \mathbf{e} = \boldsymbol{\varepsilon}^\top \mathbf{M}\boldsymbol{\varepsilon}$.

Note que os resíduos são correlacionados:

$$\text{Var}(\mathbf{e}|\mathbf{X}) = \text{E}(\mathbf{e}\mathbf{e}^\top | \mathbf{X}) = \text{E}(\mathbf{M}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top \mathbf{M} | \mathbf{X}) = \mathbf{M}\text{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top | \mathbf{X})\mathbf{M} = \mathbf{M}\sigma^2 \mathbf{I}\mathbf{M} = \sigma^2 \mathbf{M}.$$

Note que \mathbf{M} (simétrica e idempotente) é uma matriz positiva definida.

Agora, estamos interessados em

$$E(\mathbf{e}^\top \mathbf{e} | \mathbf{X}) = E(\boldsymbol{\varepsilon}^\top \mathbf{M} \boldsymbol{\varepsilon} | \mathbf{X}).$$

A quantidade escalar $\boldsymbol{\varepsilon}^\top \mathbf{M} \boldsymbol{\varepsilon}$ é uma matriz 1×1 , logo é igual ao seu traço.

Segue-se que

$$\begin{aligned} E(\text{tr}(\boldsymbol{\varepsilon}^\top \mathbf{M} \boldsymbol{\varepsilon}) | \mathbf{X}) &= E(\text{tr}(\mathbf{M} \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top) | \mathbf{X}) = \text{tr}(E(\mathbf{M} \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top | \mathbf{X})) = \text{tr}(\mathbf{M} E(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top | \mathbf{X})) \\ &= \text{tr}(\mathbf{M} \sigma^2 \mathbf{I}) = \sigma^2 \text{tr}(\mathbf{M}). \end{aligned}$$

O traço de \mathbf{M} é

$$\text{tr}(\mathbf{I}_n - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \text{tr}(\mathbf{I}_n) - \text{tr}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}) = \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{I}_p) = n - p.$$

Portanto, $E(\mathbf{e}^\top \mathbf{e} | \mathbf{X}) = (n - p)\sigma^2$.

Logo, o estimador natural de σ^2 é tendencioso. Um estimador não tendencioso de σ^2 é

$$s^2 = \frac{\mathbf{e}^\top \mathbf{e}}{n - p}.$$

Este estimador é não tendencioso incondicionalmente também,

$$E(s^2) = E_X(E(s^2|\mathbf{X})) = E_X(\sigma^2) = \sigma^2.$$

Podemos estimar $\text{Var}(\mathbf{b}|\mathbf{X})$ por $\widehat{\text{Var}}(\mathbf{b}|\mathbf{X}) = s^2(\mathbf{X}^\top \mathbf{X})^{-1}$.

A raiz quadrada do k -ésimo elemento da diagonal de $s^2(\mathbf{X}^\top \mathbf{X})^{-1}$ é o erro padrão do estimador b_k .

Normalidade (HP.6)

Se tivermos a hipótese $\varepsilon \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I})$, então

$$\begin{aligned}(\mathbf{b}|\mathbf{X}) &\sim \mathcal{N}_p(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}), \quad \text{e} \\(b_k|\mathbf{X}) &\sim \mathcal{N}(\beta_k, \sigma^2(\mathbf{X}^\top \mathbf{X})_{kk}^{-1}).\end{aligned}$$

Lembre-se que $\mathbf{b} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \boldsymbol{\beta} + \mathbf{A}\boldsymbol{\varepsilon}$ com $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$, ou seja, sob **HP.6** \mathbf{b} é uma combinação linear de normais. Portanto, \mathbf{b} (condicional a \mathbf{X}) tem distribuição normal.

Estimação por intervalo

Começaremos utilizando todas as hipóteses do modelo de regressão linear, inclusive **HP.6**.

Dado $\varepsilon \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I})$, temos que

$$\begin{aligned}(\mathbf{b}|\mathbf{X}) &\sim \mathcal{N}_p(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}), \quad \text{e} \\(b_k|\mathbf{X}) &\sim \mathcal{N}(\beta_k, \sigma^2(\mathbf{X}^\top \mathbf{X})_{kk}^{-1}).\end{aligned}$$

Definamos s_{kk} o k -ésimo elemento da diagonal de $(\mathbf{X}^\top \mathbf{X})^{-1}$.

Daí, temos

$$Z_k = \frac{b_k - \beta_k}{\sqrt{\sigma^2 s_{kk}}} \sim \mathcal{N}(0, 1).$$

Mas σ^2 é desconhecido.

Resultado 1

Os autovalores de uma matriz idempotente são compostos de 0 (zeros) ou 1 (uns).

Resultado 2

O posto de uma matriz simétrica e idempotente é igual ao seu traço (um número inteiro não negativo).

Resultado 3

A única matriz idempotente não singular é a matriz identidade.

Resultado 4

Se $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ e \mathbf{C} é uma matriz tal que $\mathbf{C}^\top \mathbf{C} = \mathbf{I}$, então $\mathbf{C}^\top \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Resultado 5

Se $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ e \mathbf{A} é uma matriz simétrica e idempotente, então $\mathbf{x}^\top \mathbf{A} \mathbf{x} \sim \chi_m^2$ sendo m o número de raízes unitárias de \mathbf{A} (que é igual ao posto de \mathbf{A}).

Resultado 6

Se $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\mathbf{L}\mathbf{x}$ uma função linear e $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ uma forma quadrática (\mathbf{A} simétrica e idempotente), então $\mathbf{L}\mathbf{x}$ e $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ são independentes se $\mathbf{L}\mathbf{A} = \mathbf{0}$.

A quantidade

$$\frac{(n-p)S^2}{\sigma^2} = \frac{\mathbf{e}^\top \mathbf{e}}{\sigma^2} = \left(\frac{\boldsymbol{\varepsilon}}{\sigma}\right)^\top \mathbf{M} \left(\frac{\boldsymbol{\varepsilon}}{\sigma}\right)$$

é uma forma quadrática simétrica e idempotente do vetor normal $(\boldsymbol{\varepsilon}/\sigma)$ com $\text{tr}(\mathbf{M}) = n - p$ e a forma linear dada por

$$\left(\frac{\mathbf{b} - \boldsymbol{\beta}}{\sigma}\right) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \left(\frac{\boldsymbol{\varepsilon}}{\sigma}\right).$$

Temos que $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{M} = \mathbf{0}$ pois $\mathbf{M}\mathbf{X} = \mathbf{0}$.

Teorema (independência de \mathbf{b} e s^2)

Se $\boldsymbol{\varepsilon}$ é um vetor com distribuição normal, então o estimador \mathbf{b} de mínimos quadrados dos coeficientes é estatisticamente independente do vetor de resíduos \mathbf{e} e, portanto, todas as funções de \mathbf{e} , incluindo s^2 .

Destes resultados, temos que

$$\frac{b_k - \beta_k}{\sqrt{s^2 s_{kk}}} \sim t_{n-p}.$$

Consequentemente,

$$\text{IC}_{100(1-\alpha)\%}(\beta_k) = (b_k - t_{n-p, 1-\alpha/2} \sqrt{s^2 s_{kk}}, b_k + t_{n-p, 1-\alpha/2} \sqrt{s^2 s_{kk}}).$$

Combinação linear dos coeficientes

- Seja \mathbf{w} um vetor $p \times 1$ de constantes.
- Estamos interessados na combinação linear $\alpha = \mathbf{w}^\top \boldsymbol{\beta}$.
- Estimamos α com $a = \mathbf{w}^\top \mathbf{b}$ sendo \mathbf{b} o estimador de mínimos quadrados.
- Como \mathbf{b} tem distribuição normal, temos

$$a \sim \mathcal{N}(\mathbf{w}^\top \boldsymbol{\beta}, \sigma^2 \mathbf{w}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{w}).$$

- Ainda estimamos σ^2 com s^2 .
- Estimamos a variância de a com $s_a^2 = s^2 \mathbf{w}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{w}$.
- Segue-se que $(a - \alpha)/s_a \sim t_{n-p}$ (pois a e s_a^2 são independentes).
- Daí, podemos construir intervalos de confiança para α .
- Podemos testar somas ou diferenças entre os coeficientes.

Predição

- Suponha que desejamos prever o valor de y_* associado ao vetor de regressores \mathbf{x}_* .
- Este valor seria $y_* = \mathbf{x}_*^\top \boldsymbol{\beta} + \varepsilon_*$.
- Segue do teorema de Gauss-Markov que $\hat{y}_* = \mathbf{x}_*^\top \mathbf{b}$ é o estimador não tendencioso linear de variância mínima de $E(y_* | \mathbf{x}_*)$.
- O erro de previsão é $e_* = y_* - \hat{y}_* = (\boldsymbol{\beta} - \mathbf{b})^\top \mathbf{x}_* + \varepsilon_*$.
- A variância de predição para este estimador é

$$\text{Var}(e_* | \mathbf{X}, \mathbf{x}_*) = \sigma^2 + \text{Var}((\boldsymbol{\beta} - \mathbf{b})^\top \mathbf{x}_* | \mathbf{X}, \mathbf{x}_*) = \sigma^2(1 + \mathbf{x}_*^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_*).$$

- Podemos estimar esta variância ao substituir σ^2 pelo estimador s^2 .
- Assim, um intervalo de predição de y_* é dada por

$$\text{IC}_{100(1-\alpha)\%}(y_*) = (\hat{y}_* - t_{\alpha/2; n-p} \times \text{ep}(e_*), \hat{y}_* + t_{\alpha/2; n-p} \times \text{ep}(e_*)).$$

Multicolinearidade

Quando os regressores são altamente correlacionados, temos que

- pequenas mudanças nos dados podem produzir grandes variações na estimativas dos parâmetros;
- coeficientes podem ter erros padrões grandes e baixos níveis de significância mesmo que eles sejam conjuntamente significativos e o R^2 para a regressão seja elevado; e
- coeficientes podem ter um sinal equivocado ou magnitudes implausíveis.

Seja \mathbf{X} a matriz de variáveis regressoras com uma constante e $p - 1$ variáveis medidas como desvios de suas médias, $\mathbf{x}_k = \mathbf{z}_k - \bar{z}_k \mathbf{1}$.

Seja $\mathbf{X}_{(k)}$ a matriz com todas as variáveis exceto \mathbf{x}_k . Então, o k -ésimo elemento de $(\mathbf{X}^\top \mathbf{X})^{-1}$ é dado por

$$\begin{aligned} (\mathbf{x}_k^\top \mathbf{M}_{(k)} \mathbf{x}_k)^{-1} &= \left[\mathbf{x}_k^\top \mathbf{x}_k - \mathbf{x}_k^\top \mathbf{X}_{(k)} (\mathbf{X}_{(k)}^\top \mathbf{X}_{(k)})^{-1} \mathbf{X}_{(k)}^\top \mathbf{x}_k \right]^{-1} \\ &= \left[\mathbf{x}_k^\top \mathbf{x}_k \left(1 - \frac{\mathbf{x}_k^\top \mathbf{X}_{(k)} (\mathbf{X}_{(k)}^\top \mathbf{X}_{(k)})^{-1} \mathbf{X}_{(k)}^\top \mathbf{x}_k}{\mathbf{x}_k^\top \mathbf{x}_k} \right) \right]^{-1} \\ &= \frac{1}{(1 - R_k^2) s_{kk}}, \end{aligned}$$

sendo R_k^2 o coeficiente de determinação da regressão de x_k sobre todas as outras variáveis do modelo e $s_{kk} = \sum_{i=1}^n (x_{ik} - \bar{x}_k)^2$.

A variância do estimador de mínimos quadrados do k -ésimo coeficiente é σ^2 vezes a razão acima,

$$\text{Var}(b_k | \mathbf{X}) = \frac{\sigma^2}{(1 - R_k^2) \sum_{i=1}^n (x_{ik} - \bar{x}_k)^2}.$$

Com as outras coisas iguais,

- quanto maior a correlação de x_k com as outras variáveis, maior será a variância de b_k , devido a multicolinearidade;
- quanto maior a variância de x_k , menor será a variância de b_k ; e
- quanto melhor o ajuste geral da regressão, menor será a variância devido a σ^2 .

Dados reais nunca serão ortogonais ($R_k^2 = 0$). Então, a multicolinearidade estará sempre presente.

Quando a multicolinearidade é um problema?

Podemos utilizar o **número de condição** de $\mathbf{X}^T \mathbf{X}$: a raiz quadrada da razão da maior raiz característica de $\mathbf{X}^T \mathbf{X}$ (sendo as colunas padronizadas para tamanho unitário) sobre a menor. Valores maiores que 20 sugerem uma indicação de problema.

Possíveis caminhos são

- aumentar o número de observações;
- retirar as variáveis mais problemáticas do modelo; e
- utilizar componentes principais (pode haver dificuldade de interpretação).

Exemplo - dados de Longley

Tabela 12: Correlação entre as variáveis.

	Emprego	Preço	PIB	Militar	Ano
Emprego	1,0000	0,9709	0,9836	0,4573	0,9713
Preço	0,9709	1,0000	0,9916	0,4647	0,9911
PIB	0,9836	0,9916	1,0000	0,4464	0,9953
Militar	0,4573	0,4647	0,4464	1,0000	0,4172
Ano	0,9713	0,9911	0,9953	0,4172	1,0000

Ver arquivo exemplo_09.r

Tabela 13: Ajuste do modelo com 15 observações.

	Estimativa	Erro padrão	Estat. <i>t</i>	Valor <i>p</i>
Constante	1459415,1	714182,9	2,04348	0,06825
Ano	-721,756	369,985	-1,95077	0,07965
Preço	-181,123	135,525	-1,33646	0,21101
PIB	0,09107	0,02026	4,49478	0,00115
Militar	-0,07494	0,26113	-0,28698	0,77999

Tabela 14: Ajuste do modelo com 16 observações.

	Estimativa	Erro padrão	Estat. <i>t</i>	Valor <i>p</i>
Constante	1169087,5	835902,4	1,39859	0,18949
Ano	-576,464	433,487	-1,32983	0,21049
Preço	-19,7681	138,893	-0,14233	0,88940
PIB	0,06439	0,01995	3,22746	0,00805
Militar	-0,01015	0,30857	-0,03288	0,97436

Testes de hipóteses

Estamos interessados em implicações testáveis. Por exemplo, o modelo apresenta algumas restrições testáveis.

Suponha o seguinte modelo de investimento (I_t):

$$M_1 : \ln I_t = \beta_1 + \beta_2 \iota_t + \beta_3 \Delta p_t + \beta_4 \ln y_t + \beta_5 t + \varepsilon_t,$$

sendo

- ι_t : taxa de juros nominal;
- Δp_t : taxa de inflação; e
- $\ln y_t$: logaritmo da produção real.

Entretanto, uma teoria alternativa afirma que “investidores se preocupam com a taxa de juros real”.

$$\ln I_t = \beta_1 + \beta_2(\iota_t - \Delta p_t) + \beta_3 \Delta p_t + \beta_4 \ln y_t + \beta_5 t + \varepsilon_t.$$

Não temos como testar a hipótese, pois ι_t e Δp_t estão em ambos os modelos.

Considere agora a alternativa em que investidores se preocupam *somente* com a taxa de juros real

$$\mathcal{M}_2 : \ln I_t = \beta_1 + \beta_2(\iota_t - \Delta p_t) + \beta_4 \ln y_t + \beta_5 t + \varepsilon_t.$$

Agora, em relação ao modelo \mathcal{M}_1 , temos a restrição $\beta_2 + \beta_3 = 0$.

Os modelos em \mathcal{M}_1 e \mathcal{M}_2 são encaixados.

Se os investidores não se preocuparem com a inflação (modelo \mathcal{M}_3), então $\beta_3 = 0$. Neste caso, o modelo \mathcal{M}_3 também seria encaixado com \mathcal{M}_1 .

Então, suponha os seguintes modelos:

Modelo \mathcal{M}_3 : Investidores se preocupam somente com a taxa de juros nominal $(\beta_1, \beta_2, 0, \beta_4, \beta_5)$.

Modelo \mathcal{M}_4 : Investidores se preocupam somente com a inflação $(\beta_1, 0, \beta_3, \beta_4, \beta_5)$.

Estes dois modelos (\mathcal{M}_3 e \mathcal{M}_4) não são encaixados. Eles têm o mesmo número de parâmetros, mas diferentes covariáveis.

A seguir trataremos de modelos encaixados.

O modelo de regressão linear

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

Restrições lineares:

$$\begin{aligned} r_{11}\beta_1 + r_{12}\beta_2 + \cdots + r_{1p}\beta_p &= q_1 \\ r_{21}\beta_1 + r_{22}\beta_2 + \cdots + r_{2p}\beta_p &= q_2 \\ &\vdots \\ r_{J1}\beta_1 + r_{J2}\beta_2 + \cdots + r_{Jp}\beta_p &= q_J. \end{aligned}$$

Podemos organizar em notação matricial $\mathbf{R}\boldsymbol{\beta} = \mathbf{q}$.

A matriz \mathbf{R} tem J linhas e p restrições, então $J < p$. Estamos eliminando a possibilidade de $J = p$ para evitar solução degenerada. As linhas de \mathbf{R} devem ser linearmente independentes.

A restrição $\mathbf{R}\boldsymbol{\beta} = \mathbf{q}$ impõe J restrições em p parâmetros livres. Então, existem $p - J$ parâmetros livres.

Abordagem para testes de hipóteses

Utilizaremos a hipótese de normalidade **HP.6**.

Queremos testar:

$$H_0 : \mathbf{R}\boldsymbol{\beta} - \mathbf{q} = \mathbf{0}$$

$$H_1 : \mathbf{R}\boldsymbol{\beta} - \mathbf{q} \neq \mathbf{0}.$$

Exemplos:

1. $\beta_j = 0$, $\mathbf{R} = [0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]$, $\mathbf{q} = 0$;
2. $\beta_k = \beta_j$, $\mathbf{R} = [0 \ \dots \ 0 \ 1 \ 0 \ \dots \ -1 \ 0 \ \dots \ 0]$, $\mathbf{q} = 0$;
3. $\beta_2 + \beta_3 + \beta_4 = 1$, $\mathbf{R} = [0 \ 1 \ 1 \ 1 \ 0 \ \dots \ 0]$, $\mathbf{q} = 1$;
4. $\beta_2 + \beta_3 = 1$, $\beta_4 + \beta_6 = 0$, $\beta_5 + \beta_6 = 0$,

$$\mathbf{R} = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix} \quad \text{e} \quad \mathbf{q} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix};$$

5. $\beta_2 = \beta_3 = \dots = \beta_p = 0$, $\mathbf{R} = [\mathbf{0} : \mathbf{I}_{p-1}]$, $\mathbf{q} = \mathbf{0}$.

Dado \mathbf{b} , estamos interessados no vetor de discrepâncias $\mathbf{m} = \mathbf{R}\mathbf{b} - \mathbf{q}$.

Em termos estatísticos, \mathbf{m} é significativamente diferente de $\mathbf{0}$?

Temos que

$$\mathbf{b} \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}) \Rightarrow \mathbf{m} \sim \mathcal{N}(\mathbf{R}\boldsymbol{\beta} - \mathbf{q}, \sigma^2 \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top).$$

Sob H_0 , $\mathbf{m} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)$ e $\mathbf{R}\boldsymbol{\beta} = \mathbf{q}$ tal que
 $(1/\sigma)[\mathbf{R}\mathbf{b} - \mathbf{q}] = (1/\sigma)\mathbf{R}(\mathbf{b} - \boldsymbol{\beta}) = \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\boldsymbol{\varepsilon}/\sigma) = \mathbf{D}(\boldsymbol{\varepsilon}/\sigma)$.

Podemos testar a H_0 com o critério de Wald. Condicional a \mathbf{X} , temos

$$\begin{aligned} W &= \mathbf{m}^\top [\text{Var}(\mathbf{m}|\mathbf{X})]^{-1} \mathbf{m} \\ &= \frac{(\mathbf{R}\mathbf{b} - \mathbf{q})^\top [\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top]^{-1} (\mathbf{R}\mathbf{b} - \mathbf{q})}{\sigma^2} \\ &= \left(\frac{\boldsymbol{\varepsilon}}{\sigma}\right)^\top \mathbf{D}^\top [\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top]^{-1} \mathbf{D} \left(\frac{\boldsymbol{\varepsilon}}{\sigma}\right), \end{aligned}$$

uma forma quadrática simétrica e idempotente em $(\boldsymbol{\varepsilon}/\sigma)$ com

$$\text{tr} \left(\mathbf{D}^\top [\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top]^{-1} \mathbf{D} \right) = \text{tr} \left([\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top]^{-1} \mathbf{D}\mathbf{D}^\top \right) = \text{tr}(\mathbf{I}_J) = J.$$

Logo, $W \sim \chi_J^2$.

Mas a estatística W não pode ser utilizada porque desconhecemos σ^2 .

Podemos utilizar s^2 (regressão) para estimar σ^2 .

Resultado 7

Dado $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, se $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ e $\mathbf{x}^\top \mathbf{B} \mathbf{x}$ são duas formas quadráticas simétricas e idempotentes em \mathbf{x} , então $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ e $\mathbf{x}^\top \mathbf{B} \mathbf{x}$ são independentes se $\mathbf{A} \mathbf{B} = \mathbf{0}$.

Note que dado $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, temos que $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ e $\mathbf{x}^\top \mathbf{B} \mathbf{x}$ têm distribuição qui-quadrado.

Sabemos que \mathbf{b} e s^2 são independentes. Isto implica que W e s^2 são independentes. Logo,

$$\begin{aligned} F &= \frac{W/J}{(n-p)s^2/[\sigma^2(n-p)]} \\ &= \frac{(\mathbf{R}\mathbf{b} - \mathbf{q})^\top [\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top]^{-1} (\mathbf{R}\mathbf{b} - \mathbf{q})}{Js^2} \sim \mathcal{F}_{J, n-p}. \end{aligned}$$

Teste de hipóteses com uma restrição

Agora, considere uma abordagem alternativa para testar a hipótese

$$H_0 : r_1\beta_1 + r_2\beta_2 + \cdots + r_p\beta_p = \mathbf{r}^\top \boldsymbol{\beta} = q.$$

Sabemos que $\hat{q} = \mathbf{r}^\top \mathbf{b}$ é o estimador de $q = \mathbf{r}^\top \boldsymbol{\beta}$ e que

$$\hat{q} = \mathbf{r}^\top \mathbf{b} \sim \mathcal{N}(\mathbf{r}^\top \boldsymbol{\beta}, \sigma^2 \mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r}).$$

Logo,

$$t = \frac{\hat{q} - q}{\text{ep}(\hat{q})} \sim t_{n-p}, \quad \text{sendo} \quad \text{ep}(\hat{q}) = s \sqrt{\mathbf{r}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{r}}.$$

Note que $t^2 \sim \mathcal{F}_{1, n-p}$. Portanto, os testes são equivalentes.

Exemplo 12 - investimento restrito

Analizaremos dados trimestrais dos EUA entre jan/1950 e dez/2000. O modelo é dado por

$$\ln I_t = \beta_1 + \beta_2 \iota_t + \beta_3 \Delta p_t + \beta_4 \ln y_t + \beta_5 t + \varepsilon_t.$$

Primeiro queremos testar $H_0 : \beta_2 + \beta_3 = 0$ (somente taxa de juros real).

- I_t é o investimento real (Realinvs);
- ι_t é a taxa de juros (Tbilrate);
- Δp_t é a inflação medida por variação no logaritmo do índice de preço ao consumidor (CPI_U); e
- $\ln y_t$ é o logaritmo do PIB (Realgdp).

Uma observação é “perdida” ao calcular Δp_t .

Ver arquivo exemplo_10.r

Tabela 15: Resumo do ajuste do modelo irrestrito.

	Estimativa	Erro padrão	Estat. t	Valor p
Constante	-9,134092	1,366459	-6,684	$< 10^{-10}$
Taxa de juros	-0,008598	0,003196	-2,691	0,00774
Inflação	0,003306	0,002337	1,415	0,15872
ln PIB	1,930156	0,183272	10,532	$< 10^{-15}$
Tempo	-0,005659	0,001488	-3,803	0,00019

Além disso, $\hat{\sigma} = 0,08618$ e $R^2 = 0,9798$.

Ainda, $F = 2.395$ com 4 e 198 graus de liberdade, e valor $p < 10^{-15}$.

Com os resultados do modelo irrestrito, temos

- $\hat{q} = b_2 + b_3 = -0,008598 + 0,003306 = -0,00529,$
- $ep(\hat{q}) = \sqrt{ep(b_2)^2 + ep(b_3)^2 + 2cov(b_2, b_3)} =$
 $\sqrt{0,00320^2 + 0,00234^2 + 2 * (-3,717e - 06)} = 0,002878, e$
- $t = \hat{q}/ep(\hat{q}) = -1,838.$
- Como $|t| < t_{0,975;198} = 1,972,$ a hipótese nula não é rejeitada.
- Podemos estimar o modelo supondo $\beta_2 = -\beta_3.$ (Próxima página.)

Ver arquivo exemplo_10.r

Tabela 16: Resumo do ajuste do modelo restrito.

	Estimativa	Erro padrão	Estat. t	Valor p
Constante	-7,90716	1,20063	-6,59	$< 10^{-10}$
Taxa de juros real	-0,00443	0,00227	-1,95	0,0526
ln PIB	1,76406	0,16056	10,99	$< 10^{-15}$
Tempo	-0,00440	0,00133	-3,31	0,0011

Além disso, $\hat{\sigma} = 0,0867$ e $R^2 = 0,979$.

Ainda, $F = 3.150$ com 3 e 198 graus de liberdade, e valor $p < 10^{-15}$.

Agora, queremos testar (conjuntamente) a seguinte hipótese nula:

- $\beta_2 + \beta_3 = 0$ (taxa de juros real),
- $\beta_4 = 1$ (a propensão marginal a investir é 1), e
- $\beta_5 = 0$ (não existe tendência temporal).

Ainda utilizando o ajuste no modelo irrestrito, temos

$$\mathbf{R} = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{q} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{m} = \mathbf{R}\mathbf{b} - \mathbf{q} = \begin{bmatrix} -0,005292 \\ 0,930156 \\ -0,005659 \end{bmatrix} \quad e$$

$$\widehat{\text{Var}}(\mathbf{m}|\mathbf{X}) = \begin{bmatrix} 8,238e-06 & -0,0002586 & 1,956e-06 \\ -2,586e-04 & 0,0335887 & -2,718e-04 \\ 1,956e-06 & -0,0002718 & 2,214e-06 \end{bmatrix}.$$

$$\text{Assim, } F = \frac{\mathbf{m}^\top [\widehat{\text{Var}}(\mathbf{m}|\mathbf{X})]^{-1} \mathbf{m}}{J} = 109,8, \text{ sendo } J = 3.$$

Como $F > F_{3; 198; 0,95} = 2,65$, a hipótese nula é rejeitada (valor $p \approx 0$).

Predição para investimento

Queremos prever o primeiro trimestre de 2001 ao nível $\alpha = 0,05$.

- Taxa de juros: 4,48 (90-day T-Bill);
- Inflação: 528,0 (CPI_U);
- PIB: 9316,8 (realGDP); e
- Tempo: 204.

Ver arquivo exemplo_10.r

O intervalo de predição é dado por:

$$IC_{95\%}(\hat{y}_{*,204}) = (7,159; 7,503).$$

A taxa anual de investimento real no primeiro trimestre de 2001 foi 1.721.

O logaritmo é 7,4507. Assim, a predição contém este valor.

Estimadores de mínimos quadrados restritos

- Sabemos que \mathbf{b} é escolhido para minimizar $\mathbf{e}^\top \mathbf{e}$.
- Sendo $R^2 = 1 - \frac{\mathbf{e}^\top \mathbf{e}}{\mathbf{y}^\top \tilde{\mathbf{M}} \mathbf{y}}$ e $\mathbf{y}^\top \tilde{\mathbf{M}} \mathbf{y}$ é uma constante que não envolve \mathbf{b} , então \mathbf{b} é escolhido para maximizar R^2 .
- O estimador de mínimos quadrados restrito é obtido pela solução de

$$\min_{\mathbf{b}} S(\mathbf{b}) = \min_{\mathbf{b}} (\mathbf{y} - \mathbf{X}\mathbf{b})^\top (\mathbf{y} - \mathbf{X}\mathbf{b}) \quad \text{sujeito a} \quad \mathbf{R}\mathbf{b} = \mathbf{q}.$$

- Uma função Lagrangeana para este problema pode ser escrita como

$$L(\mathbf{b}_*, \lambda_*) = (\mathbf{y} - \mathbf{X}\mathbf{b}_*)^\top (\mathbf{y} - \mathbf{X}\mathbf{b}_*) + 2\lambda_*^\top (\mathbf{R}\mathbf{b}_* - \mathbf{q}).$$

- Como λ_* é irrestrito, utilizamos $2\lambda_*$ como restrição.
- A solução de \mathbf{b}_* e λ_* satisfazem as condições necessárias:

$$\begin{aligned} \frac{\partial L(\mathbf{b}_*, \lambda_*)}{\partial \mathbf{b}_*} &= -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{b}_*) + 2\mathbf{R}^\top \lambda = \mathbf{0} \\ \frac{\partial L(\mathbf{b}_*, \lambda_*)}{\partial \lambda_*} &= 2(\mathbf{R}\mathbf{b}_* - \mathbf{q}) = \mathbf{0}. \end{aligned}$$

$$\begin{bmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{R}^T \\ \mathbf{R} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{b}_* \\ \lambda_* \end{bmatrix} = \begin{bmatrix} \mathbf{X}^T \mathbf{y} \\ \mathbf{q} \end{bmatrix} \quad \text{ou} \quad \mathbf{A} \mathbf{d}_* = \mathbf{v}.$$

Se a matriz \mathbf{A} for não singular, então $\mathbf{d}_* = \mathbf{A}^{-1} \mathbf{v}$.

Se $\mathbf{X}^T \mathbf{X}$ for não singular, então uma solução explícita para \mathbf{b}_* e λ_* pode ser obtida usando a fórmula da inversa particionada:

$$\mathbf{b}_* = \mathbf{b} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \left[\mathbf{R} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \right]^{-1} (\mathbf{R} \mathbf{b} - \mathbf{q}) = \mathbf{b} - \mathbf{C} \mathbf{m} \quad \text{e}$$

$$\lambda_* = \left[\mathbf{R} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \right]^{-1} (\mathbf{R} \mathbf{b} - \mathbf{q}).$$

Além disso,

$$\text{Var}(\mathbf{b}_* | \mathbf{X}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} - \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \left[\mathbf{R} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \right]^{-1} \mathbf{R} (\mathbf{X}^T \mathbf{X})^{-1}.$$

(A prova está na próxima página.)

Logo, $\text{Var}(\mathbf{b}_* | \mathbf{X})$ é a variância de \mathbf{b} reduzida pelas restrições.

Prova:

$$\begin{aligned}\text{Var}(\mathbf{b}_*|\mathbf{X}) &= \text{Var}(\mathbf{b}|\mathbf{X}) + \text{Var}(\mathbf{C}\mathbf{m}|\mathbf{X}) - 2\text{Cov}(\mathbf{b}, \mathbf{C}\mathbf{m}|\mathbf{X}) \\ &= \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{C}\text{Var}(\mathbf{R}\mathbf{b} - \mathbf{q}|\mathbf{X})\mathbf{C}^\top - 2\text{Cov}(\mathbf{b}, \mathbf{R}\mathbf{b} - \mathbf{q}|\mathbf{X})\mathbf{C}^\top \\ &= \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{C}\mathbf{R}\sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \mathbf{C}^\top - 2\text{Cov}(\mathbf{b}, \mathbf{b}|\mathbf{X})\mathbf{R}^\top \mathbf{C}^\top \\ &= \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} + \sigma^2 \mathbf{C}\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \mathbf{C}^\top - 2\sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \mathbf{C}^\top.\end{aligned}$$

Note que

$$\begin{aligned}\mathbf{C}\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top &= (\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \left[\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \right]^{-1} \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \\ &= (\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top.\end{aligned}$$

Daí, temos que

$$\begin{aligned}\text{Var}(\mathbf{b}_*|\mathbf{X}) &= \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} + \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \mathbf{C}^\top - 2\sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \mathbf{C}^\top \\ &= \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} - \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \mathbf{C}^\top \\ &= \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} - \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \left[\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{R}^\top \right]^{-1} \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1}.\end{aligned}$$

A solução restrita de mínimos quadrados não pode ser melhor que a irrestrita:

$$\begin{aligned} \mathbf{e}_* &= \mathbf{y} - \mathbf{X}\mathbf{b}_* = \mathbf{y} - \mathbf{X}\mathbf{b} + \mathbf{X}\mathbf{b} - \mathbf{X}\mathbf{b}_* = \mathbf{e} - \mathbf{X}(\mathbf{b}_* - \mathbf{b}) \Rightarrow \\ \mathbf{e}_*^\top \mathbf{e}_* &= \mathbf{e}^\top \mathbf{e} + (\mathbf{b}_* - \mathbf{b})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{b}_* - \mathbf{b}) \geq \mathbf{e}^\top \mathbf{e} \Rightarrow \\ R_*^2 &\leq R^2. \end{aligned}$$

A perda de ajuste é dada por

$$\mathbf{e}_*^\top \mathbf{e}_* - \mathbf{e}^\top \mathbf{e} = (\mathbf{R}\mathbf{b} - \mathbf{q})^\top \left[\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \right]^{-1} (\mathbf{R}\mathbf{b} - \mathbf{q}).$$

Esta é a expressão no numerador do teste \mathcal{F} :

$$\begin{aligned} F &= \frac{(\mathbf{R}\mathbf{b} - \mathbf{q})^\top \left[\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \right]^{-1} (\mathbf{R}\mathbf{b} - \mathbf{q})}{J s^2} = \frac{(\mathbf{e}_*^\top \mathbf{e}_* - \mathbf{e}^\top \mathbf{e})/J}{\mathbf{e}^\top \mathbf{e}/(n-p)} \\ &= \frac{(R^2 - R_*^2)/J}{(1 - R^2)/(n-p)} \sim \mathcal{F}_{J, n-p}. \end{aligned}$$

Se a restrição for $\beta_2 = \dots = \beta_p = 0$ (exceto intercepto), então temos $p - 1$ restrições e o teste \mathcal{F} usual de significância da regressão.

Exemplo - função de produção

Função de produção de Cobb-Douglas para a indústria de metal.

Regressão de mínimos quadrados no logaritmo da produção (valor adicionado) sobre uma constante, logaritmo do trabalho e produção de capital.

Uma generalização do modelo de Cobb-Douglas é o modelo **translog**:

$$\begin{aligned} \ln y &= \beta_1 + \beta_2 \ln TR + \beta_3 \ln CP \\ &+ \beta_4 \left[\frac{1}{2} (\ln TR)^2 \right] + \beta_5 \left[\frac{1}{2} (\ln CP)^2 \right] + \beta_6 \ln TR \times \ln CP + \varepsilon, \end{aligned}$$

sendo TR a variável trabalho e CP a variável capital.

O modelo de Cobb-Douglas é obtido com $\beta_4 = \beta_5 = \beta_6 = 0$.

Ver arquivo exemplo_11.r

Tabela 17: Resumo do ajuste do modelo translog.

	Estimativa	Erro padrão	Estat. t	Valor p
Constante	0,9442	2,9108	0,324	0,7489
ln TR	3,6136	1,5481	2,334	0,0296
ln CP	-1,8931	1,0163	-1,863	0,0765
$(\ln TR)^2/2$	-0,9641	0,7074	-1,363	0,1874
$(\ln CP)^2/2$	0,0853	0,2926	0,291	0,7735
$(\ln TR)(\ln CP)$	0,3124	0,4389	0,712	0,4845

Além disso, $\hat{\sigma} = 0,1799$ e $R^2 = 0,9549$.

Ainda, $F = 88,85$ com 5 e 21 graus de liberdade, e valor $p < 10^{-12}$.

Note o valor da estatística do teste \mathcal{F} e seu valor p . Os demais resultados e outros testes estão descritos no arquivo **R**.

Forma funcional e mudança de estrutura

Analisaremos maneiras de mudar a estrutura da variável dependente - no contexto de regressão - utilizando covariáveis.

Dados dicotômicos:

- os dados dicotômicos (binário, *dummy*) são muito utilizados em análise de regressão;
- em geral, utiliza-se estas variáveis em conjunto com outras contínuas.

Exemplo - rendimentos

- $T = 428$ participantes do mercado de trabalho formal;
- Uma equação de (semilog) ganhos:

$$\ln(\text{rendimentos}) = \beta_1 + \beta_2 \text{idade} + \beta_3 \text{idade}^2 + \beta_4 \text{educação} + \beta_5 \text{filhos} + \varepsilon;$$

- rendimentos: salário por hora vezes horas trabalhadas;
- educação: anos de escola; e
- filhos: variável dicotômica indicando (1) a presença de filhos menores de 18 anos.

Ver arquivo exemplo_12.r

Tabela 18: Resumo do ajuste do modelo de rendimentos.

Parâmetro	Estimativa	Erro padrão	Estat. t	Valor p
β_1	3,2401	1,7674	1,833	0,0675
β_2	0,2006	0,0839	2,392	0,0172
β_3	-0,0023	0,0010	-2,345	0,0195
β_3	0,0675	0,0253	2,672	0,0078
β_5	-0,3512	0,1475	-2,380	0,0177

Além disso, $\hat{\sigma} = 1,19$ e $R^2 = 0,041$.

O valor $-0,3512$ é considerado como um efeito extremamente grande (semilog).

Note que $\epsilon = \frac{\partial y/y}{\partial x/x} = \frac{\partial \ln y}{\partial \text{filhos}} \times \text{filhos} = \beta_k \text{filhos}$ (é um modelo semilog).

$1 - \exp(-0,3512) = 0,2962 \simeq 1/3$ (mulheres com filhos recebem quase 1/3 menos).

Possíveis aplicações de variáveis dicotômicas

- Estudo do efeito de tratamentos em algum tipo de resposta.
- Diferenças de sexo em relação ao salário na iniciativa privada.
- Modelo envolvendo uma variável dicotômica:

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \delta d_i + \varepsilon_i.$$

- É prática comum incluir variáveis dicotômicas na regressão para modelar algo somente em uma observação.

Uma variável dicotômica que toma valor 1 (um) somente para uma observação tem o efeito de excluir esta observação do cálculo dos coeficientes de mínimos quadrados e a variância do estimador (mas não o R^2).

- Diversas categorias (por exemplo, efeitos sazonais):

$$C_t = \beta_1 + \beta_2 x_t + \delta_1 D_{t1} + \delta_2 D_{t2} + \delta_3 D_{t3} + \varepsilon_t,$$

sendo x_t os ganhos.

- Uma formulação alternativa é dada por

$$C_t = \beta x_t + \delta_1 D_{t1} + \delta_2 D_{t2} + \delta_3 D_{t3} + \delta_4 D_{t4} + \varepsilon_t.$$

- Diversos grupos de variáveis dicotômicas.

$$y_{it} = \alpha + \beta x_{it} + \delta_i + \theta_t + \varepsilon_{it}.$$

- Precisamos retirar uma dicotômica de cada grupo.

Exemplo - companhias aéreas

O estudo trata da eficiência na produção de serviços de companhias aéreas.

Temos observações de 6 firmas durante 15 anos.

O modelo é dado por

$$\begin{aligned} \ln C_{i,t} &= \beta_1 + \beta_2 \ln Q_{i,t} + \beta_3 (\ln Q_{i,t})^2 + \beta_4 \ln P_{i,t} + \beta_5 \text{Cargas}_{i,t} \\ &+ \sum_{t=1}^{14} \theta_t D_{i,t} + \sum_{i=1}^5 \delta_i F_{i,t} + \varepsilon_{i,t} \end{aligned}$$

Ver arquivo exemplo_13.r

Forma funcional e mudança de estrutura (continuação)

Em várias aplicações, variáveis dicotômicas são utilizadas para fatores qualitativos somente.

Suponha o interesse no seguinte modelo:

$$\text{renda} = \beta_1 + \beta_2 \text{idade} + \text{efeito de educação} + \varepsilon.$$

Ensino médio (EM), graduação (GR), mestrado (MS) e doutorado (DS).

Uma maneira **insatisfatória** de modelar é considerar a variável E igual a 0 para o primeiro grupo, 1 para o segundo, 2 para o terceiro e 3 para o quarto.

$$\text{renda} = \beta_1 + \beta_2 \text{idade} + \beta_3 E + \varepsilon.$$

Um modelo mais flexível seria um com variáveis dicotômicas:

$$\text{renda} = \beta_1 + \beta_2 \text{idade} + \delta_1 \text{GR} + \delta_2 \text{MS} + \delta_3 \text{DS} + \varepsilon.$$

Ensino Médio: $E(\text{renda}|\text{idade}, \text{EM}) = \beta_1 + \beta_2 \text{idade},$

Graduação: $E(\text{renda}|\text{idade}, \text{GR}) = \beta_1 + \beta_2 \text{idade} + \delta_1,$

Mestrado: $E(\text{renda}|\text{idade}, \text{MS}) = \beta_1 + \beta_2 \text{idade} + \delta_2,$

Doutorado: $E(\text{renda}|\text{idade}, \text{DS}) = \beta_1 + \beta_2 \text{idade} + \delta_3.$

Regressão *spline* (função definida segmentarmente por polinômios)

A função que desejamos estimar é

$$\begin{aligned} E(\text{renda}|\text{idade}) &= \alpha^0 + \beta^0 \text{idade}, & \text{se } \text{idade} < 18, \\ & \alpha^1 + \beta^1 \text{idade}, & \text{se } 18 \leq \text{idade} < 22, \\ & \alpha^2 + \beta^2 \text{idade}, & \text{se } \text{idade} \geq 22. \end{aligned}$$

Os valores 18 e 22 são chamados de nós. Seja $d_1 = 1$, se $\text{idade} \geq t_1^*$ e $d_2 = 1$, se $\text{idade} \geq t_2^*$, sendo $t_1^* = 18$ e $t_2^* = 22$. Temos a combinação das equações dada por

$$\text{renda} = \beta_1 + \beta_2 \text{idade} + \gamma_1 d_1 + \delta_1 d_1 \text{idade} + \gamma_2 d_2 + \delta_2 d_2 \text{idade} + \varepsilon.$$

Os coeficientes angulares dos três segmentos são β_2 , $\beta_2 + \delta_1$ e $\beta_2 + \delta_1 + \delta_2$.

Para fazer com que a função seja contínua, precisamos impor que

$$\begin{aligned} \beta_1 + \beta_2 t_1^* &= (\beta_1 + \gamma_1) + (\beta_2 + \delta_1) t_1^* & \text{e} \\ (\beta_1 + \gamma_1) + (\beta_2 + \delta_1) t_2^* &= (\beta_1 + \gamma_1 + \gamma_2) + (\beta_2 + \delta_1 + \delta_2) t_2^*. \end{aligned}$$

Isto implica restrições lineares nos coeficientes:

$$\gamma_1 + \delta_1 t_1^* = 0 \quad \text{and} \quad \gamma_2 + \delta_2 t_2^* = 0.$$

Assim,

$$\text{renda} = \beta_1 + \beta_2 \text{idade} + \delta_1 d_1 (\text{idade} - t_1^*) + \delta_2 d_2 (\text{idade} - t_2^*) + \varepsilon.$$

Os estimadores de mínimos quadrados restritos são obtidos pela regressão múltipla:

$$x_1 = \text{idade},$$

$$x_2 = \text{idade} - 18, \quad \text{se } \text{idade} \geq 18 \quad \text{e } 0 \quad \text{caso contrário},$$

$$x_3 = \text{idade} - 22, \quad \text{se } \text{idade} \geq 22 \quad \text{e } 0 \quad \text{caso contrário}.$$

Podemos testar a hipótese nula $H_0 : \delta_1 = \delta_2 = 0$ que implica mesmo coeficiente angular.

Não linearidade nas variáveis

Seja $\mathbf{z} = \{z_1, z_2, \dots, z_L\}$ um conjunto de L variáveis independentes.

Seja f_1, f_2, \dots, f_p um conjunto de p funções linearmente independentes de \mathbf{z} .

Seja $g(y)$ uma função observável de y .

$$\begin{aligned}g(y) &= \beta_1 f_1(\mathbf{z}) + \beta_2 f_2(\mathbf{z}) + \dots + \beta_p f_p(\mathbf{z}) + \varepsilon \\ &= \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon \\ &= \mathbf{x}^T \boldsymbol{\beta} + \varepsilon.\end{aligned}$$

Um modelo comum é o log-linear

$$\ln y = \ln \alpha + \sum_k \beta_k \ln z_k + \varepsilon = \beta_1 + \sum_k \beta_k x_k + \varepsilon.$$

No modelo log-linear,

- Mudanças são medidas em termos proporcionais ou porcentagem.
- β_k mede a mudança percentual em y associada com a mudança de 1 por cento em x_k .

Termos de interação

Um modelo relacionando distância de frenagem D com velocidade S e humidade da pista W ,

$$D = \beta_1 + \beta_2 S + \beta_3 W + \beta_4 SW + \varepsilon.$$

Neste modelo,

$$\frac{\partial E(D|S, W)}{\partial S} = \beta_2 + \beta_4 W.$$

Isto implica que o efeito marginal da velocidade na distância de frenagem cresce quando a pista fica mais molhada.

$$\text{Var} \left(\frac{\partial \hat{E}(D|S, W)}{\partial S} \right) = \text{Var}(\hat{\beta}_2) + W^2 \text{Var}(\hat{\beta}_4) + 2W \text{Cov}(\hat{\beta}_2, \hat{\beta}_4).$$

Podemos tentar identificar não linearidades:

- Análise de resíduos;
- Regressão linear por partes; e
- Regressão polinomial.

Exemplo - função de produção de elasticidade constante de substituição

O modelo é dado por

$$\ln y = \ln \gamma - \frac{\nu}{\rho} \ln(\delta K^{-\rho} + (1 - \delta)L^{-\rho}) + \varepsilon$$

A aproximação de Taylor em torno do ponto $\rho = 0$ é

$$\begin{aligned} \ln y &= \ln \gamma + \nu \delta \ln K + \nu(1 - \delta) \ln L + \rho \nu \delta(1 - \delta) \left\{ -\frac{1}{2} [\ln K - \ln L]^2 \right\} + \varepsilon^* \\ &= \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \varepsilon^*, \end{aligned}$$

sendo $x_1 = 1$, $x_2 = \ln K$, $x_3 = \ln L$ e $x_4 = -(1/2) \ln^2(K/L)$.

As transformações são dadas por

$$\begin{aligned} \beta_1 &= \ln \gamma & \gamma &= \exp(\beta_1) \\ \beta_2 &= \nu \delta & \delta &= \beta_2 / (\beta_2 + \beta_3) \\ \beta_3 &= \nu(1 - \delta) & \nu &= \beta_2 + \beta_3 \\ \beta_4 &= \rho \nu \delta(1 - \delta) & \rho &= \beta_4 (\beta_2 + \beta_3) / (\beta_2 \beta_3). \end{aligned}$$

Estimativas de β_1 , β_2 , β_3 e β_4 podem ser obtidas por mínimos quadrados.

Para estimar a matriz de covariância de $\theta = (\gamma, \delta, \nu, \rho)$ utilizamos o **método delta** (ver próxima página).

$$\mathbf{C} = \frac{\partial \theta}{\partial \beta^T} = \begin{bmatrix} \exp(\beta_1) & 0 & 0 & 0 \\ 0 & \frac{\beta_3}{(\beta_2 + \beta_3)^2} & -\frac{\beta_2}{(\beta_2 + \beta_3)^2} & 0 \\ 0 & 1 & 1 & 0 \\ 0 & -\frac{\beta_3\beta_4}{\beta_2^2\beta_3} & -\frac{\beta_2\beta_4}{\beta_2\beta_3^2} & \frac{\beta_2 + \beta_3}{\beta_2\beta_3} \end{bmatrix}.$$

A matriz de covariâncias de θ é, então, estimada por

$$\hat{\mathbf{C}} \left[s^2 (\mathbf{X}^T \mathbf{X})^{-1} \right] \hat{\mathbf{C}}^T.$$

Vale ressaltar que nem todos os modelos são intrinsecamente lineares.

Método delta

Caso univariado: Seja X_n uma sequência de variáveis aleatórias tal que

$$X_n \xrightarrow{P} \theta \quad \text{e} \quad \sqrt{n}(X_n - \theta) \xrightarrow{D} \mathcal{N}(0, \sigma^2)$$

então

$$\sqrt{n}(g(X_n) - g(\theta)) \xrightarrow{D} \mathcal{N}(0, \sigma^2 [g'(\theta)]^2)$$

desde que $g'(\theta)$ exista e seja diferente de zero.

Caso multivariado: Seja \mathbf{X}_n uma sequência de vetores aleatórios tal que

$$\mathbf{X}_n \xrightarrow{P} \boldsymbol{\theta} \quad \text{e} \quad \sqrt{n}(\mathbf{X}_n - \boldsymbol{\theta}) \xrightarrow{D} \mathcal{N}_d(\mathbf{0}, \boldsymbol{\Sigma})$$

então

$$\sqrt{n}(\mathbf{h}(\mathbf{X}_n) - \mathbf{h}(\boldsymbol{\theta})) \xrightarrow{D} \mathcal{N}_k(\mathbf{0}, \nabla \mathbf{h}(\boldsymbol{\theta})^\top \boldsymbol{\Sigma} \nabla \mathbf{h}(\boldsymbol{\theta}))$$

desde que $\nabla \mathbf{h}(\boldsymbol{\theta})$ exista e seja bem definido.

Modelo de regressão generalizado e heterocedasticidade

Trataremos de modelos que não satisfazem a hipótese **HP.4** de homocedasticidade e correlação nula.

O modelo de regressão linear generalizado é dado por

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ E(\boldsymbol{\varepsilon}|\mathbf{X}) &= \mathbf{0} \\ E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top|\mathbf{X}) &= \sigma^2\boldsymbol{\Omega} = \boldsymbol{\Sigma}, \end{aligned}$$

em que $\boldsymbol{\Sigma}$ é uma matriz positiva definida.

No modelo (somente) heterocedástico:

$$\sigma^2\boldsymbol{\Omega} = \sigma^2 \begin{bmatrix} \omega_1 & 0 & \cdots & 0 \\ 0 & \omega_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \omega_n \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix}.$$

No modelo (somente) correlacionado:

$$\sigma^2 \mathbf{\Omega} = \sigma^2 \begin{bmatrix} 1 & \rho_1 & \cdots & \rho_{n-1} \\ \rho_1 & 1 & \cdots & \rho_{n-2} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{n-1} & \rho_{n-2} & \cdots & 1 \end{bmatrix}.$$

Em geral, aparece em modelos de séries temporais (modelo AR(1)).

Conjuntos de dados de painéis, consistindo de observações por seções cruzadas em vários pontos no tempo, podem exibir heterocedasticidade e correlação não nula.

Estimação por mínimos quadrados

Para o modelo com distúrbios esféricos,

$$E(\boldsymbol{\varepsilon}|\mathbf{X}) = \mathbf{0} \quad \text{e} \quad E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top|\mathbf{X}) = \sigma^2\mathbf{I},$$

o estimador de mínimos quadrados ordinários (OLS) é

$$\mathbf{b} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}.$$

- o melhor estimador linear não tendencioso (BLUE);
- consistente;
- assintoticamente normal; e
- se os erros forem normais, então também é eficiente assintoticamente.

Teorema

Se os regressores e os distúrbios forem não correlacionados, então o estimador de mínimos quadrados para o modelo de regressão linear generalizado é não tendencioso. Condicional a \mathbf{X} , a variância amostral do estimador de mínimos quadrados é

$$\begin{aligned}\text{Var}(\mathbf{b}|\mathbf{X}) &= \text{E}((\mathbf{b} - \beta)(\mathbf{b} - \beta)^\top | \mathbf{X}) \\ &= \text{E}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \varepsilon \varepsilon^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} | \mathbf{X}) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \sigma^2 \boldsymbol{\Omega} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \frac{\sigma^2}{n} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \right)^{-1} \left(\frac{1}{n} \mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X} \right) \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \right)^{-1}\end{aligned}$$

Se ε for normalmente distribuído, então

$$\mathbf{b} \sim \mathcal{N}(\beta, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X}) (\mathbf{X}^\top \mathbf{X})^{-1}).$$

- Não podemos mais estimar $\text{Var}(\mathbf{b}|\mathbf{X})$ com $s^2 (\mathbf{X}^\top \mathbf{X})^{-1}$; e
- Os testes t e F já não são mais apropriados.

Propriedades assintóticas de mínimos quadrados

Teorema

Se $\mathbf{Q} = \text{plim}(\mathbf{X}^\top \mathbf{X}/n)$ e $\text{plim}(\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X}/n)$ são matrizes positivas definidas e finitas, então \mathbf{b} é consistente para $\boldsymbol{\beta}$. Sob estas hipóteses, $\text{plim} \mathbf{b} = \boldsymbol{\beta}$.

A condição no teorema acima depende de \mathbf{X} e de $\boldsymbol{\Omega}$.

Teorema

Se os regressores forem suficientemente bem comportados e os termos fora da diagonal da matriz $\boldsymbol{\Omega}$ diminuírem suficientemente rápido, então o estimador de mínimos quadrados é assintoticamente normal com média $\boldsymbol{\beta}$ e matriz de covariâncias

$$\text{Var}(\mathbf{b}) = \frac{\sigma^2}{n} \mathbf{Q}^{-1} \text{plim} \left(\frac{1}{n} \mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X} \right) \mathbf{Q}^{-1}.$$

Estimação robusta das matrizes de covariâncias assintóticas

Se $\sigma^2\Omega$ fosse conhecido, então o estimador da matriz de covariância assintótica de \mathbf{b} seria

$$\mathbf{V}_{OLS} = \frac{1}{n} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \right)^{-1} \left(\frac{1}{n} \mathbf{X}^\top \left[\sigma^2 \Omega \right] \mathbf{X} \right) \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \right)^{-1}.$$

Nosso problema é $\sigma^2\Omega$. Assumiremos que $\text{tr}(\Omega) = n$ e $\Sigma = \sigma^2\Omega$.

Em princípio, parece que temos $n(n+1)/2$ parâmetros para estimar em Σ .

Na verdade, precisamos estimar $p(p+1)/2$ parâmetros da matriz

$$\text{plim} \mathbf{Q}_* = \text{plim} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \sigma_{ij} \mathbf{x}_i \mathbf{x}_j^\top.$$

Voltaremos a este ponto mais à frente.

Estimação eficiente por mínimos quadrados generalizados

Temos que Ω é uma matriz positiva definida, logo

$$\Omega = \mathbf{C}\mathbf{\Lambda}\mathbf{C}^T,$$

sendo as colunas de \mathbf{C} os vetores característicos de Ω e $\mathbf{\Lambda}$ é uma matriz diagonal com as raízes características de Ω .

Seja $\mathbf{T} = \mathbf{C}\mathbf{\Lambda}^{1/2}$, então $\Omega = \mathbf{T}\mathbf{T}^T$.

Seja $\mathbf{P}^T = \mathbf{C}\mathbf{\Lambda}^{-1/2}$, então $\Omega^{-1} = \mathbf{P}^T\mathbf{P}$.

Segue-se que

$$\mathbf{P}\mathbf{y} = \mathbf{P}\mathbf{X}\boldsymbol{\beta} + \mathbf{P}\boldsymbol{\varepsilon} \quad \text{ou} \quad \mathbf{y}_* = \mathbf{X}_*\boldsymbol{\beta} + \boldsymbol{\varepsilon}_* \Rightarrow \text{E}(\boldsymbol{\varepsilon}_*\boldsymbol{\varepsilon}_*^T | \mathbf{X}_*) = \mathbf{P}\sigma^2\mathbf{\Omega}\mathbf{P}^T = \sigma^2\mathbf{I}.$$

Então, o modelo de regressão linear clássico se aplica a este modelo transformado.

Por hipótese, Ω é conhecido, então \mathbf{y}_* e \mathbf{X}_* são os dados observados. Logo,

$$\hat{\beta} = (\mathbf{X}_*^T \mathbf{X}_*)^{-1} \mathbf{X}_*^T \mathbf{y}_* = (\mathbf{X}^T \mathbf{P}^T \mathbf{P} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{P}^T \mathbf{P} \mathbf{y} = (\mathbf{X}^T \Omega^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Omega^{-1} \mathbf{y}$$

é um estimador eficiente de β . Este é o estimador de mínimos quadrados generalizados (GLS).

Teorema

1. Se $E(\varepsilon_* | \mathbf{X}_*) = \mathbf{0} \Leftrightarrow E(\mathbf{P}\varepsilon | \mathbf{P}\mathbf{X}) = \mathbf{0} \Leftrightarrow E(\varepsilon | \mathbf{X}) = \mathbf{0}$, então $E(\hat{\beta} | \mathbf{X}_*) = \beta$.
2. O estimador GLS é consistente se $\text{plim}(1/n) \mathbf{X}_*^T \mathbf{X}_* = \mathbf{Q}_*$, sendo \mathbf{Q}_* uma matriz positiva definida. Logo, $\text{plim}[(1/n) \mathbf{X}^T \Omega^{-1} \mathbf{X}]^{-1} = \mathbf{Q}_*^{-1}$.
3. O estimador GLS é normalmente distribuído com média β e a variância amostral

$$\text{Var}(\hat{\beta} | \mathbf{X}_*) = \sigma^2 (\mathbf{X}_*^T \mathbf{X}_*)^{-1} = \sigma^2 (\mathbf{X}^T \Omega \mathbf{X})^{-1}.$$

4. O estimador GLS $\hat{\beta}$ é o estimador linear não tendencioso de variância mínima no modelo de regressão generalizado.

Para testar J restrições lineares, $\mathbf{R}\beta = \mathbf{q}$, a estatística apropriada é

$$F_{J, n-p} = \frac{(\mathbf{R}\hat{\beta} - \mathbf{q})^\top [\mathbf{R}\hat{\sigma}^2(\mathbf{X}_*^\top \mathbf{X}_*)^{-1} \mathbf{R}^\top]^{-1} (\mathbf{R}\hat{\beta} - \mathbf{q})}{J} = \frac{(\hat{\varepsilon}_c^\top \hat{\varepsilon}_c - \hat{\varepsilon}^\top \hat{\varepsilon})/J}{\hat{\sigma}^2},$$

sendo $\hat{\varepsilon} = \mathbf{y}_* - \mathbf{X}_* \hat{\beta}$ e $\hat{\sigma}^2 = \frac{\hat{\varepsilon}^\top \hat{\varepsilon}}{n-p} = \frac{(\mathbf{y} - \mathbf{X}\hat{\beta})^\top \Omega^{-1} (\mathbf{y} - \mathbf{X}\hat{\beta})}{n-p}$.

Os resíduos do GLS restrito, $\hat{\varepsilon}_c = \mathbf{y}_* - \mathbf{X}_* \hat{\beta}_c$, são baseados em

$$\hat{\beta}_c = \hat{\beta} - [\mathbf{X}^\top \Omega^{-1} \mathbf{X}]^{-1} [\mathbf{R}(\mathbf{X}^\top \Omega^{-1} \mathbf{X})^{-1} \mathbf{R}^\top]^{-1} (\mathbf{R}\hat{\beta} - \mathbf{q}).$$

Não existe nenhuma contrapartida precisa para o R^2 no modelo de regressão generalizado.

Se Ω for desconhecido, então não é possível utilizar o estimador GLS.

Precisamos impor restrições em Ω para diminuir o número efetivo de parâmetros a estimar.

Heterocedasticidade

Distúrbios da regressão cujas variâncias não são constantes para diferentes observações são heterocedásticos.

Em modelos de regressão heterocedásticos,

$$\text{Var}(\varepsilon_i | \mathbf{X}) = \sigma_i^2 = \sigma^2 \omega_i, \quad i = 1, 2, \dots, n \quad \text{ou}$$

$$\mathbf{E}(\varepsilon \varepsilon^\top) = \sigma^2 \mathbf{\Omega} = \sigma^2 \begin{bmatrix} \omega_1 & 0 & \cdots & 0 \\ 0 & \omega_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \omega_n \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix},$$

sendo $\text{tr}(\mathbf{\Omega}) = \sum_{i=1}^n \omega_i = n$ para evitar problemas de escala.

Ineficiência de mínimos quadrados

- \mathbf{b} , o estimador OLS, é ineficiente relativo ao estimador GLS.
- Quanto maior a dispersão de ω_i para diferentes observações, maior a eficiência do GLS sobre OLS.
- Suponha que a forma da heterocedasticidade seja desconhecida. Então, o estimador GLS não pode ser usado.
- Se usarmos \mathbf{b} , então a matriz de covariância $\sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$ é inapropriada.
- A matriz $\sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}(\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X})(\mathbf{X}^\top \mathbf{X})^{-1}$ é apropriada.
- Temos, usualmente,

$$s^2 = \frac{\mathbf{e}^\top \mathbf{e}}{n-p} = \frac{\boldsymbol{\varepsilon}^\top \mathbf{M} \boldsymbol{\varepsilon}}{n-p} = \frac{\boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}}{n-p} - \frac{\boldsymbol{\varepsilon}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}}{n-p},$$

pois $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$. Isto implica que

$$\mathbb{E} \left(\frac{\boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}}{n-p} \mid \mathbf{X} \right) = \frac{\text{tr}(\mathbb{E}(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top \mid \mathbf{X}))}{n-p} = \frac{n\sigma^2}{n-p} \quad \text{e}$$

$$\begin{aligned}
E\left(\frac{\boldsymbol{\varepsilon}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}}{n-p} \mid \mathbf{X}\right) &= \frac{\text{tr}(E[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top \mathbf{X}])}{n-p} \\
&= \frac{\text{tr}(\sigma^2(\mathbf{X}^\top \mathbf{X}/n)^{-1}(\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X}/n))}{n-p} \\
&= \frac{\sigma^2}{n-p} \text{tr}\left(\left(\frac{\mathbf{X}^\top \mathbf{X}}{n}\right)^{-1} \mathbf{Q}_n^*\right),
\end{aligned}$$

sendo $\mathbf{Q}_n^* = \frac{\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X}}{n} = \frac{1}{n} \sum_{i=1}^n \omega_i \mathbf{x}_i \mathbf{x}_i^\top$.

Se \mathbf{b} é consistente, então $\lim_{n \rightarrow \infty} E(s^2) = \sigma^2$.

Se a heterocedasticidade não é correlacionada com as variáveis no modelo, então pelo menos em amostras grandes, os cálculos de mínimos quadrados, apesar de não serem a maneira ótima de utilizar os dados, não levarão a “erros grosseiros”.

Observações:

- É verdade que o OLS é ineficiente.
- Se por hipótese da análise - que a heterocedasticidade não é relacionada as variáveis do modelo - é incorreta, então os erros padrões convencionais podem estar longe dos valores apropriados.

Estamos interessados em um estimador para \mathbf{Q}_* = $\frac{1}{n} \sum_{i=1}^n \sigma_i^2 \mathbf{x}_i \mathbf{x}_i^\top$.

Pode ser mostrado sob condições bem gerais que o estimador

$$\mathbf{S}_0 = \frac{1}{n} \sum_{i=1}^n e_i^2 \mathbf{x}_i \mathbf{x}_i^\top$$

tem $\text{plim} \mathbf{S}_0 = \text{plim} \mathbf{Q}_*$.

Assim, o estimador consistente de heterocedasticidade de **White** é

$$\widehat{\text{Var}}(\mathbf{b}|\mathbf{X}) = n(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{S}_0 (\mathbf{X}^\top \mathbf{X})^{-1}.$$

Um dos problemas dos quadrados dos resíduos de OLS é que eles tendem a subestimar os quadrados dos distúrbios verdadeiros.

Davidson e McKinnon propuseram

1. escalar o resultado final por $n/(n - p)$; e
2. usar e_i^2/m_{ii} ao invés de e_i^2 , sendo $m_{ii} = 1 - \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i$.

Exemplo - estimadores de White, e de Davidson e McKinnon

Ver arquivo exemplo_14.r

Modelaremos os gastos mensais com cartão de crédito (G):

$$G = \beta_1 + \beta_2 \text{idade} + \beta_3 \text{casa} + \beta_4 \text{renda} + \beta_5 \text{renda}^2 + \varepsilon.$$

A variável “casa” é uma binária indicando casa própria ou não.

Testando a heterocedasticidade

O teste geral de White é dado por $H_0 : \sigma_i^2 = \sigma^2$ para todo i , e $H_1 : \text{Não } H_0$.

A dificuldade é em estimar um modelo com n parâmetros.

A matriz de covariância correta para o estimador de mínimos quadrados é

$$\text{Var}(\mathbf{b}|\mathbf{X}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}(\mathbf{X}^\top \mathbf{\Omega} \mathbf{X})(\mathbf{X}^\top \mathbf{X})^{-1}$$

que pode ser estimada como anteriormente.

O estimador convencional é $\mathbf{V} = s^2(\mathbf{X}^\top \mathbf{X})^{-1}$.

Se não existir heterocedasticidade, então \mathbf{V} dará um estimador consistente de $\text{Var}(\mathbf{b}|\mathbf{X})$.

Uma versão simples e operacional do teste de White é feita ao obter nR^2 na regressão de e_i^2 sobre uma constante, todas as variáveis em \mathbf{x}_i e todos os quadrados e produtos cruzados de \mathbf{x}_i .

A estatística tem distribuição assintótica qui-quadrada com $r - 1$ graus de liberdade, sendo r o número de regressores na equação, incluindo a constante.

O teste de White não é construtivo. Se H_0 é rejeitada, então o resultado do teste não dá nenhuma indicação do que fazer a seguir.

O teste de **Breusch-Pagan** é baseado nos multiplicadores de Lagrange e tem $H_0 : \sigma_i^2 = \sigma^2 f(\alpha_0 + \boldsymbol{\alpha}^\top \mathbf{z}_i)$, sendo \mathbf{z}_i o vetor de variáveis independentes.

O modelo é homocedástico se $\boldsymbol{\alpha} = \mathbf{0}$.

O teste pode ser conduzido com uma regressão simples.

A estatística do teste é

$$LM = \frac{1}{2} \mathbf{g}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{g},$$

sendo \mathbf{g} o vetor de observações $g_i = e_i^2 / (\mathbf{e}^\top \mathbf{e} / n) - 1$.

Sob H_0 , $LM \approx \chi_r^2$, sendo r o número de variáveis em \mathbf{z} .

Este teste é sensível a hipótese de normalidade.

Koenker e Basset sugeriram utilizar um estimador robusto para a variância de ε_j^2 ,

$$V = \frac{1}{n} \sum_{i=1}^n \left[e_i^2 - \frac{\mathbf{e}^\top \mathbf{e}}{n} \right]^2.$$

Assim,

$$LM = \frac{1}{V} (\mathbf{u} - \bar{u}\mathbf{1})^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top (\mathbf{u} - \bar{u}\mathbf{1}),$$

sendo $\mathbf{u} = (e_1^2, e_2^2, \dots, e_n^2)$ e $\bar{u} = \mathbf{e}^\top \mathbf{e} / n$.

Ver arquivo [exemplo_14.r](#)

Mínimos quadrados ponderados

Tendo testado e encontrado evidência de heterocedasticidade, o próximo passo é levar isto em conta na estimação.

O estimador GLS é

$$\hat{\beta} = (\mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{y}.$$

Suponha que $\text{Var}(\varepsilon_i | \mathbf{X}) = \sigma_i^2 = \sigma^2 \omega_i$.

Então, o i -ésimo elemento da diagonal da matriz $\mathbf{\Omega}^{-1}$ é $1/\omega_i$.

O estimador GLS é obtido ao regredir

$$\mathbf{Py} = \begin{bmatrix} y_1 / \sqrt{\omega_1} \\ y_2 / \sqrt{\omega_2} \\ \vdots \\ y_n / \sqrt{\omega_n} \end{bmatrix} \quad \text{sobre} \quad \mathbf{PX} = \begin{bmatrix} \mathbf{x}_1^T / \sqrt{\omega_1} \\ \mathbf{x}_2^T / \sqrt{\omega_2} \\ \vdots \\ \mathbf{x}_n^T / \sqrt{\omega_n} \end{bmatrix}.$$

Aplicando-se o OLS ao modelo transformado, nós obtemos o estimador de mínimos quadrados ponderados (WLS),

$$\hat{\beta} = \left[\sum_{i=1}^n w_i \mathbf{x}_i \mathbf{x}_i^\top \right]^{-1} \left[\sum_{i=1}^n w_i \mathbf{x}_i y_i \right],$$

sendo $w_i = 1/\omega_i$.

Uma especificação comum é que a variância seja proporcional a um dos regressores ou de seu quadrado, por exemplo se $\sigma_i^2 = \sigma^2 x_{ik}^2$, então o modelo de regressão transformado para GLS é

$$\frac{y}{x_k} = \beta_k + \beta_1 \frac{x_1}{x_k} + \beta_2 \frac{x_2}{x_k} + \dots + \frac{\varepsilon}{x_k}.$$

O estimador WLS é consistente seja qual for os pesos utilizados, mas os pesos não podem ser correlacionados com os distúrbios.

Entretanto, a escolha dos pesos deve ser feita com cuidado para não “afetar em demasia” os erros padrões do estimador GLS.

Estimação por máxima verossilhança

A função de verossilhança é dada por

$$L(\boldsymbol{\theta}|\mathbf{y}) = f(y_1, \dots, y_n|\boldsymbol{\theta}) = \prod_{i=1}^n f(y_i|\boldsymbol{\theta})$$

para y_i , $i = 1, \dots, n$, iid condicional em $\boldsymbol{\theta}$.

O logaritmo da função de verossilhança é dada por

$$\ell(\boldsymbol{\theta}|\mathbf{y}) = \ln L(\boldsymbol{\theta}|\mathbf{y}) = \sum_{i=1}^n \ln f(y_i|\boldsymbol{\theta}).$$

Suponha que no modelo clássico de regressão os distúrbios sejam normais. Então,

$$\ell(\boldsymbol{\theta}|\mathbf{y}, \mathbf{X}) = \sum_{i=1}^n \ln f(y_i|\mathbf{x}_i, \boldsymbol{\theta}) = -\frac{1}{2} \sum_{i=1}^n \left[\ln \sigma^2 + \ln(2\pi) + \frac{1}{\sigma^2} (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right]$$

sendo \mathbf{X} a matriz $n \times p$ dos regressores com a i -ésima linha igual a \mathbf{x}_i .

Estamos interessados em obter estimativas dos parâmetros, θ .

Definição

O vetor de parâmetros θ é **identificável** (estimável) se para qualquer outro vetor de parâmetros, $\gamma \neq \theta$, para algum dado \mathbf{y} , $L(\gamma|\mathbf{y}) \neq L(\theta|\mathbf{y})$.

O princípio da máxima verossimilhança

Suponha uma amostra de 10 observações de uma distribuição Poisson: 5, 0, 1, 1, 0, 3, 2, 3, 4 e 1.

$$L(\theta|\mathbf{y}) = \frac{e^{-100\theta} \theta^{20}}{207.360}.$$

Qual valor mais provável de θ que produziu esta amostra?

A função tem um único máximo em $\theta = 2$, que é a estimativa de máxima verossimilhança.

$$\begin{aligned}\ell(\theta|\mathbf{y}) &= -n\theta + n\bar{y} \ln \theta - \sum_{i=1}^n \ln(y_i!) \\ \frac{\partial \ell(\theta|\mathbf{y})}{\partial \theta} &= -n + \frac{n\bar{y}}{\theta} = 0, \quad \Rightarrow \quad \hat{\theta}_{MV} = \bar{y}_n \\ \frac{\partial^2 \ell(\theta|\mathbf{y})}{\partial \theta^2} &= -\frac{n\bar{y}}{\theta^2} < 0, \quad \Rightarrow \quad \bar{y}_n \text{ é ponto de máximo.}\end{aligned}$$

Chamamos de **equação de verossimilhança** a expressão:

$$\frac{\partial \ell(\theta|\mathbf{dados})}{\partial \theta} = \frac{\partial \ln L(\theta|\mathbf{dados})}{\partial \theta} = \mathbf{0}.$$

Denotaremos:

- $\hat{\theta}$ como o estimador de máxima verossimilhança;
- θ_0 o valor verdadeiro do vetor de parâmetros; e
- θ um outro valor possível do vetor parâmetros.

Teorema

Sob regularidade, o EMV tem as seguintes propriedades assintóticas:

M1: Consistência: $\text{plim } \hat{\theta} = \theta_0$.

M2: Normalidade: $\hat{\theta} \approx \mathcal{N}(\theta_0, [I(\theta_0)]^{-1})$ sendo

$$I(\theta_0) = E \left(- \frac{\partial^2 \ell(\theta | \mathbf{y})}{\partial \theta \partial \theta^\top} \Big| \theta_0 \right).$$

M3: Eficiência: $\hat{\theta}$ é assintoticamente eficiente e atinge o limite inferior de Cramér-Rao.

M4: Invariância: o EMV de $\gamma = \mathbf{c}(\theta_0)$ é $\mathbf{c}(\hat{\theta})$ se $\mathbf{c}(\theta_0)$ é uma função contínua e continuamente diferenciável (transformações biunívocas).

Se $\mathbf{g}_i = \frac{\partial \ln f(y_i | \theta)}{\partial \theta}$ (a função escore) e $\mathbf{H}_i = \frac{\partial^2 \ln f(y_i | \theta)}{\partial \theta \partial \theta^\top}$ (a matriz de informação), então

D1: $E(\mathbf{g}_i(\theta) | \theta_0) = \mathbf{0}$; e

D2: $\text{Var}(\mathbf{g}_i(\theta) | \theta_0) = E(-\mathbf{H}_i(\theta) | \theta_0)$.

Teorema

$E \left(\frac{1}{n} \ell(\boldsymbol{\theta}_0) \middle| \boldsymbol{\theta}_0 \right) > E \left(\frac{1}{n} \ell(\boldsymbol{\theta}) \middle| \boldsymbol{\theta}_0 \right)$ para qualquer $\boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ (incluindo $\hat{\boldsymbol{\theta}}$).

O valor esperado da log-verossimilhança é maximizado no valor verdadeiro dos parâmetros.

Teorema (Cramér-Rao)

Sob condições de regularidade, a variância assintótica de um estimador, consistente e assintoticamente normal, do vetor de parâmetros $\boldsymbol{\theta}_0$ será sempre no mínimo maior ou igual a

$$\begin{aligned} [I(\boldsymbol{\theta}_0)]^{-1} &= \left[E \left(- \frac{\partial^2 \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}_0 \partial \boldsymbol{\theta}_0^\top} \middle| \boldsymbol{\theta}_0 \right) \right]^{-1} \\ &= \left[E \left(\left(\frac{\partial \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}_0} \right) \left(\frac{\partial \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}_0} \right)^\top \middle| \boldsymbol{\theta}_0 \right) \right]^{-1}. \end{aligned}$$

Teste da razão de verossimilhança

Seja $\hat{\theta}_U$ o EMV sob o modelo irrestrito e $\hat{\theta}_R$ sob o modelo restrito.

Sejam também \hat{L}_U e \hat{L}_R as respectivas funções de verossimilhança avaliadas nestes pontos.

Então, a razão de verossimilhança é

$$\hat{\lambda} = \frac{\hat{L}_R}{\hat{L}_U}, \quad \text{com } 0 \leq \lambda \leq 1.$$

Teorema

Sob regularidade e sob H_0 , a distribuição para amostras grandes de $-2 \ln \hat{\lambda}$ é qui-quadrada, com os graus de liberdade igual ao número de restrições impostas.

A hipótese nula é rejeitada se este valor ultrapassa o valor crítico apropriado da distribuição qui-quadrada.

Comparando modelos

Duas medidas comuns para modelos não encaixados baseados na mesma lógica são

- Critério de Informação de Akaike: $AIC = -2 \ln L(\hat{\theta}) + 2p$; e
- Critério de Informação Bayesiana/Schwarz: $BIC = -2 \ln L(\hat{\theta}) + p \ln n$,

sendo p o número de parâmetros do modelo. Menor AIC (ou BIC), melhor o modelo.

Modelo de regressão linear normal

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i.$$

$$\begin{aligned} L(\boldsymbol{\beta}, \sigma^2) &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right). \end{aligned}$$

$$\ell(\boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2}.$$

$$\begin{bmatrix} \frac{\partial \ell(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta}} \\ \frac{\partial \ell(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2} \end{bmatrix} = \begin{bmatrix} \frac{\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{\sigma^2} \\ -\frac{n}{2\sigma^2} + \frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^4} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix} \Rightarrow$$

$$\hat{\boldsymbol{\beta}}_{MV} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad \text{e} \quad \hat{\sigma}_{MV}^2 = \frac{\mathbf{e}^\top \mathbf{e}}{n}.$$

Temos,

$$\begin{bmatrix} \frac{\partial^2 \ell(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} & \frac{\partial^2 \ell(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta} \partial \sigma^2} \\ \frac{\partial^2 \ell(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2 \partial \boldsymbol{\beta}} & \frac{\partial^2 \ell(\boldsymbol{\beta}, \sigma^2)}{\partial (\sigma^2)^2} \end{bmatrix} = \begin{bmatrix} -\frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2} & -\frac{\mathbf{X}^\top \boldsymbol{\varepsilon}}{\sigma^4} \\ -\frac{\boldsymbol{\varepsilon}^\top \mathbf{X}}{\sigma^4} & \frac{n}{2\sigma^4} - \frac{\boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon}}{\sigma^6} \end{bmatrix} \Rightarrow$$

$$\left[\mathbf{I}(\boldsymbol{\beta}, \sigma^2) \right]^{-1} = \begin{bmatrix} \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} & \mathbf{0} \\ \mathbf{0} & \frac{2\sigma^4}{n} \end{bmatrix}.$$

Temos, $\sqrt{n}(\hat{\sigma}_{MV}^2 - \sigma^2) \xrightarrow{d} \mathcal{N}(0, 2\sigma^4)$.

Teste F

Se $H_0 : \mathbf{R}\boldsymbol{\beta} - \mathbf{q} = \mathbf{0}$, o teste da razão F ,

$$F = \frac{(\mathbf{R}\mathbf{b} - \mathbf{q})^\top [\mathbf{R}\mathbf{s}^2(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top](\mathbf{R}\mathbf{b} - \mathbf{q})}{J} \sim \mathcal{F}_{J, n-p},$$

para qualquer tamanho de amostra se os distúrbios forem normalmente distribuídos.

Os outros testes e suas estatísticas de teste continuam não tendo distribuição exata conhecida para amostras finitas ou pequenas.

Modelo de regressão generalizado

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, 2, \dots, n$$

$$E(\varepsilon | \mathbf{X}) = \mathbf{0}$$

$$E(\varepsilon \varepsilon^\top | \mathbf{X}) = \sigma^2 \boldsymbol{\Omega}.$$

Por hipótese, teremos $\boldsymbol{\Omega}$ como uma matriz constante conhecida.

$$\ell(\boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2} \ln |\boldsymbol{\Omega}| - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Omega}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

$$\frac{\partial \ell}{\partial \boldsymbol{\beta}} = \frac{1}{\sigma^2} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \frac{1}{\sigma^2} \mathbf{X}_*^\top (\mathbf{y}_* - \mathbf{X}_* \boldsymbol{\beta}) = \mathbf{0} \quad \boldsymbol{\Omega}^{-1} = \mathbf{P}^\top \mathbf{P}$$

$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{\sigma^4} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Omega}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \quad \mathbf{X}_* = \mathbf{P}\mathbf{X}$$

$$= -\frac{n}{2\sigma^2} + \frac{1}{\sigma^4} (\mathbf{y}_* - \mathbf{X}_* \boldsymbol{\beta})^\top (\mathbf{y}_* - \mathbf{X}_* \boldsymbol{\beta}) = 0. \quad \mathbf{y}_* = \mathbf{P}\mathbf{y}$$

$$\hat{\boldsymbol{\beta}}_{MV} = (\mathbf{X}_*^\top \mathbf{X}_*)^{-1} \mathbf{X}_*^\top \mathbf{y}_* = (\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{y} \quad \text{e}$$

$$\hat{\sigma}_{MV}^2 = \frac{1}{n} (\mathbf{y}_* - \mathbf{X}_* \hat{\boldsymbol{\beta}}_{MV})^\top (\mathbf{y}_* - \mathbf{X}_* \hat{\boldsymbol{\beta}}_{MV})$$

$$= \frac{1}{n} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{MV})^\top \boldsymbol{\Omega}^{-1} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{MV}).$$

Temos que $\hat{\sigma}_{MV}^2$ é tendencioso para σ^2 . Para obter um estimador não tendencioso precisamos multiplicá-lo pelo fator $n/(n-p)$.

Se Ω for desconhecido, então precisamos estimar $(\beta, \sigma^2, \Omega)$ simultaneamente.

Mas Ω tem $n(n+1)/2 - 1$ parâmetros. Precisamos impor restrições.

Em geral, a estimação conjunta de todos os parâmetros será complicada.

Modelo de heterocedasticidade multiplicativa

Considere um modelo de regressão com variância dada por

$$\sigma_i^2 = \sigma^2 \exp(\mathbf{w}_i^\top \boldsymbol{\alpha}) = \exp(\mathbf{z}_i^\top \boldsymbol{\gamma}),$$

sendo $\mathbf{z}_i^\top = (1, \mathbf{w}_i^\top)$ e $\boldsymbol{\gamma}^\top = (\ln \sigma^2, \boldsymbol{\alpha}^\top)$.

Neste caso, tomamos $\Omega = \text{diag}(\exp(\mathbf{z}_1^\top \boldsymbol{\gamma}), \exp(\mathbf{z}_2^\top \boldsymbol{\gamma}), \dots, \exp(\mathbf{z}_n^\top \boldsymbol{\gamma}))$.

$$\begin{aligned}
\ell(\boldsymbol{\beta}, \boldsymbol{\gamma}) &= -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^n \ln \sigma_i^2 - \frac{1}{2} \sum_{i=1}^n \frac{\varepsilon_i^2}{\sigma_i^2} \\
&= -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^n \mathbf{z}_i^\top \boldsymbol{\gamma} - \frac{1}{2} \sum_{i=1}^n \frac{\varepsilon_i^2}{\exp(\mathbf{z}_i^\top \boldsymbol{\gamma})}
\end{aligned}$$

$$\frac{\partial \ell(\boldsymbol{\beta}, \boldsymbol{\gamma})}{\partial \boldsymbol{\beta}} = \sum_{i=1}^n \mathbf{x}_i \frac{\varepsilon_i}{\exp(\mathbf{z}_i^\top \boldsymbol{\gamma})} = \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \boldsymbol{\varepsilon} = \mathbf{0}$$

$$\frac{\partial \ell(\boldsymbol{\beta}, \boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}} = \frac{1}{2} \sum_{i=1}^n \mathbf{z}_i \left(\frac{\varepsilon_i^2}{\exp(\mathbf{z}_i^\top \boldsymbol{\gamma})} - 1 \right) = \mathbf{0}.$$

$$\frac{\partial^2 \ell(\boldsymbol{\beta}, \boldsymbol{\gamma})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} = -\sum_{i=1}^n \frac{1}{\exp(\mathbf{z}_i^\top \boldsymbol{\gamma})} \mathbf{x}_i \mathbf{x}_i^\top = -\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X}$$

$$\frac{\partial^2 \ell(\boldsymbol{\beta}, \boldsymbol{\gamma})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\gamma}^\top} = -\sum_{i=1}^n \frac{\varepsilon_i}{\exp(\mathbf{z}_i^\top \boldsymbol{\gamma})} \mathbf{x}_i \mathbf{z}_i^\top$$

$$\frac{\partial^2 \ell(\boldsymbol{\beta}, \boldsymbol{\gamma})}{\partial \boldsymbol{\gamma} \partial \boldsymbol{\gamma}^\top} = -\frac{1}{2} \sum_{i=1}^n \frac{\varepsilon_i^2}{\exp(\mathbf{z}_i^\top \boldsymbol{\gamma})} \mathbf{z}_i \mathbf{z}_i^\top.$$

Temos,

$$\begin{aligned} E\left(-\frac{\partial^2 \ell(\boldsymbol{\beta}, \boldsymbol{\gamma})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\gamma}^\top}\right) &= \mathbf{0} \text{ porque } E(\varepsilon_i | \mathbf{x}_i, \mathbf{z}_i) = 0, \text{ e} \\ E\left(-\frac{\partial^2 \ell(\boldsymbol{\beta}, \boldsymbol{\gamma})}{\partial \boldsymbol{\gamma} \partial \boldsymbol{\gamma}^\top}\right) &= \frac{1}{2} \mathbf{Z}^\top \mathbf{Z} \text{ porque } E\left(\frac{\varepsilon_i^2}{\sigma_i^2} \middle| \mathbf{x}_i, \mathbf{z}_i\right) = 1. \end{aligned}$$

Seja $\boldsymbol{\delta} = (\boldsymbol{\beta}, \boldsymbol{\gamma})$. Então,

$$E\left(-\frac{\partial^2 \ell(\boldsymbol{\delta})}{\partial \boldsymbol{\delta} \partial \boldsymbol{\delta}^\top}\right) = \begin{bmatrix} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \frac{1}{2} \mathbf{Z}^\top \mathbf{Z} \end{bmatrix} = \mathbf{I}(\boldsymbol{\delta}).$$

O método do escore é um algoritmo para encontrar uma solução:

$$\boldsymbol{\delta}_{t+1} = \boldsymbol{\delta}_t + [\mathbf{I}(\boldsymbol{\delta}_t)]^{-1} \nabla \ell(\boldsymbol{\delta}_t)$$

sendo $\nabla \ell(\boldsymbol{\delta}) = \left(\frac{\partial \ell(\boldsymbol{\delta})}{\partial \boldsymbol{\beta}^\top}, \frac{\partial \ell(\boldsymbol{\delta})}{\partial \boldsymbol{\gamma}^\top} \right)^\top$.

Como $I(\delta)$ é bloco diagonal, temos

$$\begin{aligned}\beta_{t+1} &= \beta_t + (\mathbf{X}^\top \Omega_t^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Omega_t^{-1} \varepsilon_t \\ &= \beta_t + (\mathbf{X}^\top \Omega_t^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Omega_t^{-1} (\mathbf{y} - \mathbf{X}\beta_t) \\ &= (\mathbf{X}^\top \Omega_t^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Omega_t^{-1} \mathbf{y} \\ \gamma_{t+1} &= \gamma_t + \left[2(\mathbf{Z}^\top \mathbf{Z})^{-1} \right] \left[\frac{1}{2} \sum_{i=1}^n \mathbf{z}_i \left(\frac{\varepsilon_{i,t}}{\exp(\mathbf{z}_i^\top \gamma_t)} - 1 \right) \right].\end{aligned}$$

Esboço do algoritmo:

1. Estime a variância dos distúrbios σ_i^2 com $\exp(\mathbf{z}_i^\top \gamma)$.
2. Calcule β_{t+1} .
3. Calcule γ_{t+1} .
4. Calcule $\mathbf{d}_{t+1} = \|\delta_{t+1} - \delta_t\|$. Se $\mathbf{d}_{t+1} > \epsilon$, retorne ao Passo 1.

Temos também que

$$\begin{aligned}\text{Var}(\widehat{\beta}_{MV}) &= (\mathbf{X}^\top \Omega^{-1} \mathbf{X})^{-1} \\ \text{Var}(\widehat{\gamma}_{MV}) &= 2(\mathbf{Z}^\top \mathbf{Z})^{-1}.\end{aligned}$$

Estimação e inferência bayesiana

$$\begin{aligned} p(\theta|\mathbf{y}) &= \frac{p(\mathbf{y}|\theta)p(\theta)}{p(\mathbf{y})} = \frac{L(\theta|\mathbf{y})p(\theta)}{p(\mathbf{y})} \\ &\propto p(\mathbf{y}|\theta)p(\theta) = L(\theta|\mathbf{y})p(\theta), \end{aligned}$$

sendo

- $p(\theta|\mathbf{y})$ a distribuição a posteriori;
- $p(\theta)$ a distribuição a priori;
- $p(\mathbf{y}|\theta) = L(\theta|\mathbf{y})$ a função de verossimilhança; e
- $p(\mathbf{y})$ a distribuição marginal dos dados ou a verossimilhança marginal.

Temos,

$$p(\mathbf{y}) = \int p(\mathbf{y}|\theta)p(\theta)d\theta.$$

Além disso, $p(\mathbf{y}) < \infty$ para existência da distribuição a posteriori.

Análise bayesiana do modelo de regressão clássico

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

Logo,

$$L(\boldsymbol{\beta}, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}.$$

Seja $d = n - p$ os graus de liberdade e

$$\mathbf{y} - \mathbf{X}\boldsymbol{\beta} = \mathbf{y} - \mathbf{X}\mathbf{b} - \mathbf{X}(\boldsymbol{\beta} - \mathbf{b}) = \mathbf{e} - \mathbf{X}(\boldsymbol{\beta} - \mathbf{b}).$$

Então,

$$-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \left(-\frac{dS^2}{2} \right) \sigma^{-2} + \frac{1}{2} (\boldsymbol{\beta} - \mathbf{b})^\top [\sigma^{-2} \mathbf{X}^\top \mathbf{X}] (\boldsymbol{\beta} - \mathbf{b}).$$

Logo,

$$\begin{aligned} L(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) &= (2\pi)^{-d/2} (\sigma^2)^{-d/2} \exp \left\{ -\frac{d\mathbf{s}^2}{2\sigma^2} \right\} \\ &\times (2\pi)^{-p/2} (\sigma^2)^{-p/2} \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta} - \mathbf{b})^\top [\sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}]^{-1} (\boldsymbol{\beta} - \mathbf{b}) \right\} \\ &\propto \frac{(\nu \mathbf{s}^2)^{\nu-1}}{\Gamma(\nu-1)} (\sigma^2)^{-\nu} \exp \left\{ -\frac{\nu \mathbf{s}^2}{\sigma^2} \right\} \\ &\times (2\pi)^{-p/2} |\sigma^2 \boldsymbol{\Delta}|^{-1/2} \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta} - \mathbf{b})^\top [\sigma^2 \boldsymbol{\Delta}]^{-1} (\boldsymbol{\beta} - \mathbf{b}) \right\}, \end{aligned}$$

sendo $n/2 = d/2 + p/2$, $\nu = d/2$ e $\boldsymbol{\Delta} = (\mathbf{X}^\top \mathbf{X})^{-1}$.

A verossimilhança é proporcional ao produto de uma inversa gama (σ^2) com parâmetros $\delta = \nu - 1$ e $\lambda = \nu \mathbf{s}^2$, e uma normal p dimensional ($\boldsymbol{\beta} | \sigma^2$) com média \mathbf{b} e matriz de covariância $\sigma^2 \boldsymbol{\Delta}$.

Podemos considerar uma priori não informativa,

$$p(\boldsymbol{\beta}, \sigma^2) \propto \frac{1}{\sigma^2}.$$

$$\begin{aligned}
p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) &\propto L(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) p(\boldsymbol{\beta}, \sigma^2) \\
&\propto \frac{(\nu \mathbf{S}^2)^\nu}{\Gamma(\nu)} (\sigma^2)^{-(\nu+1)} \exp \left\{ -\frac{\nu \mathbf{S}^2}{\sigma^2} \right\} \\
&\quad \times (2\pi)^{-p/2} |\sigma^2 \boldsymbol{\Delta}|^{-1/2} \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta} - \mathbf{b})^\top [\sigma^2 \boldsymbol{\Delta}]^{-1} (\boldsymbol{\beta} - \mathbf{b}) \right\}.
\end{aligned}$$

Agora seja $A(\boldsymbol{\beta}) = (\boldsymbol{\beta} - \mathbf{b})^\top \boldsymbol{\Delta}^{-1} (\boldsymbol{\beta} - \mathbf{b})/2$, então

$$\begin{aligned}
p(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) &= \int p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) d\sigma^2 \propto \int_0^\infty (\sigma^2)^{-(n/2+1)} \exp \left\{ -\frac{A(\boldsymbol{\beta}) + \nu \mathbf{S}^2}{\sigma^2} \right\} d\sigma^2 \\
&= \frac{\Gamma(n/2)}{[A(\boldsymbol{\beta}) + \nu \mathbf{S}^2]^{n/2}} \propto \Gamma(n/2) \left[1 + \frac{1}{2\nu \mathbf{S}^2} A(\boldsymbol{\beta}) \right]^{-((n-p)+p)/2} \\
&= \Gamma(n/2) \left[1 + \frac{1}{2\nu} (\boldsymbol{\beta} - \mathbf{b})^\top [\mathbf{S}^2 (\mathbf{X}^\top \mathbf{X})^{-1}]^{-1} (\boldsymbol{\beta} - \mathbf{b}) \right]^{-((n-p)+p)/2}.
\end{aligned}$$

Assim,

$$(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) \sim t_{n-p} \left(\mathbf{b}, \mathbf{S}^2 (\mathbf{X}^\top \mathbf{X})^{-1} \right) \quad (p \text{ dimensional}).$$

É fácil mostrar que

$$(\sigma^2 | \mathbf{y}, \mathbf{X}) \sim \mathcal{IG}(\nu, \nu s^2).$$

Logo,

$$\begin{aligned} E(\beta | \mathbf{y}, \mathbf{X}) &= \mathbf{b}, \quad \text{para } n - p > 1, \\ E(\sigma^2 | \mathbf{y}, \mathbf{X}) &= \frac{\nu}{\nu - 1} s^2, \quad \text{para } \nu > 1 \text{ (} n - p > 2 \text{)}. \end{aligned}$$

Podemos considerar também uma priori informativa. Por exemplo,

$$(\beta | \sigma^2) \sim \mathcal{N}_p(\mathbf{b}_0, \sigma^2 \mathbf{B}_0) \quad \text{e} \quad \sigma^2 \sim \mathcal{IG}(m, m\sigma_0^2) \quad (\text{priori normal gama}).$$

$$\begin{aligned} p(\beta, \sigma^2 | \mathbf{y}, \mathbf{X}) &\propto (\sigma^2)^{-(\nu+m+1)} \exp \left\{ -\frac{1}{\sigma^2} (m\sigma_0^2 + \nu s^2) \right\} \\ &\times \left| \boldsymbol{\Sigma} \right|^{-1/2} \exp \left\{ -\frac{1}{2} (\beta - \mathbf{a})^\top \boldsymbol{\Sigma}^{-1} (\beta - \mathbf{a}) \right\}, \end{aligned}$$

sendo

$$\begin{aligned} \boldsymbol{\Sigma} &= \left[\sigma^{-2} \mathbf{B}_0^{-1} + \sigma^{-2} \mathbf{X}^\top \mathbf{X} \right]^{-1} = \sigma^2 \left[\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X} \right]^{-1} \quad \text{e} \\ \mathbf{a} &= \boldsymbol{\Sigma} \left(\sigma^{-2} \mathbf{B}_0^{-1} \mathbf{b}_0 + \sigma^{-2} \mathbf{X}^\top \mathbf{X} \mathbf{b} \right). \end{aligned}$$

Assim,

$$\begin{aligned}(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) &\sim t_{n-p}(\mathbf{a}, \boldsymbol{\Sigma}^*) \text{ e} \\(\sigma^2|\mathbf{y}, \mathbf{X}) &\sim \mathcal{IG}(\nu + m, m\sigma_0^2 + \nu s^2).\end{aligned}$$

Exercício

Encontre a expressão exata de $\boldsymbol{\Sigma}^*$.

Em problemas mais complicados, teremos que recorrer a métodos numéricos, como a amostragem de Gibbs.

Para utilizar a amostragem de Gibbs no modelo de regressão com priori não informativa, temos

$$\begin{aligned}(\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}) &\sim \mathcal{N}_p(\mathbf{b}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}) \text{ e} \\(\sigma^2|\boldsymbol{\beta}, \mathbf{y}, \mathbf{X}) &\sim \mathcal{IG}(n/2, C(\boldsymbol{\beta})),\end{aligned}$$

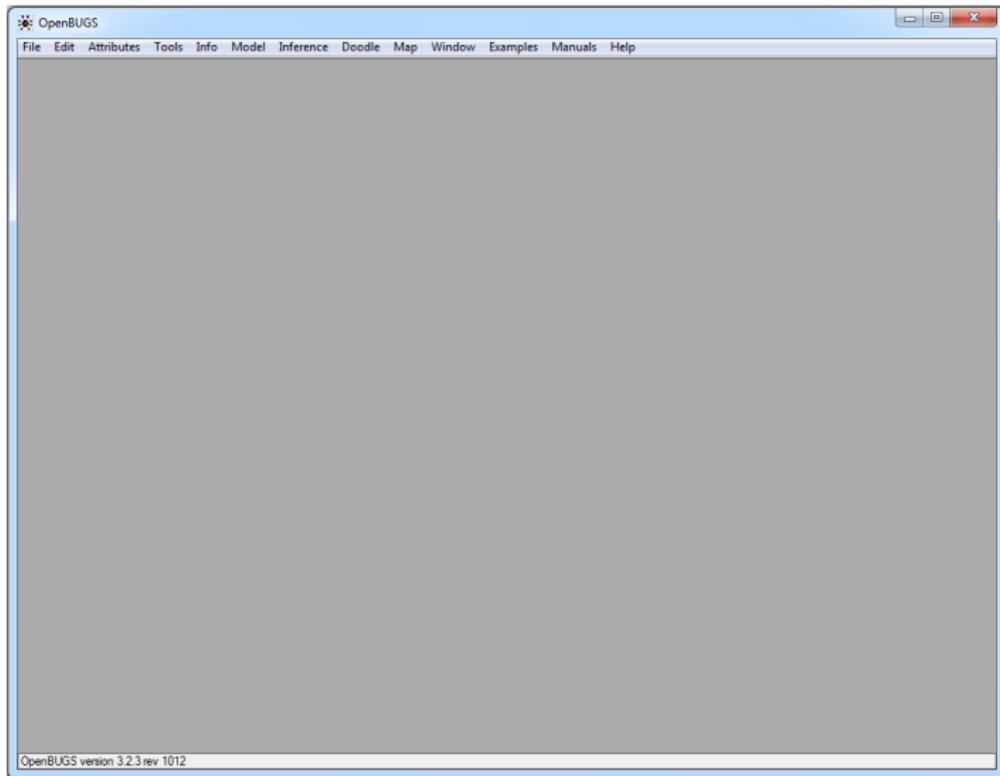
sendo $C(\boldsymbol{\beta}) = \nu s^2 + (\boldsymbol{\beta} - \mathbf{b})^\top \mathbf{X}^\top \mathbf{X}(\boldsymbol{\beta} - \mathbf{b})$.

O OpenBUGS

- Bastante amigável desde o ponto de vista do usuário;
- **OpenBUGS** disponível no sítio [the-bugs-project-openbugs](#);
- É uma ferramenta ideal para a modelagem;
- Permite mudar a especificação da distribuição a priori de forma simples;
- Em algumas aplicações a convergência pode ser lenta.

Antes, utilizamos o **WinBUGS** disponível [the-bugs-project-winbugs](#);

O entorno do OpenBUGS



Exemplo

Considere o modelo de regressão linear múltiplo dado por:

$$y_t = \beta_0 + \beta_1 x_{1t} + \beta_2 x_{2t} + \cdots + \beta_p x_{pt} + \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, \sigma^2), \quad t = 1, 2, \dots, n.$$

A distribuição a posteriori é dada por:

$$f(\beta, \sigma^2 | \mathbf{y}) \propto (\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \right\} f(\beta, \sigma^2),$$

em que $f(\beta, \sigma^2)$ é a distribuição a priori, $\mathbf{y} = (y_1, \dots, y_n)^\top$, $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$, $\mathbf{x}_i = (x_{i1}, \dots, x_{in})^\top$, e $\beta = (\beta_0, \beta_1, \dots, \beta_p)^\top$.

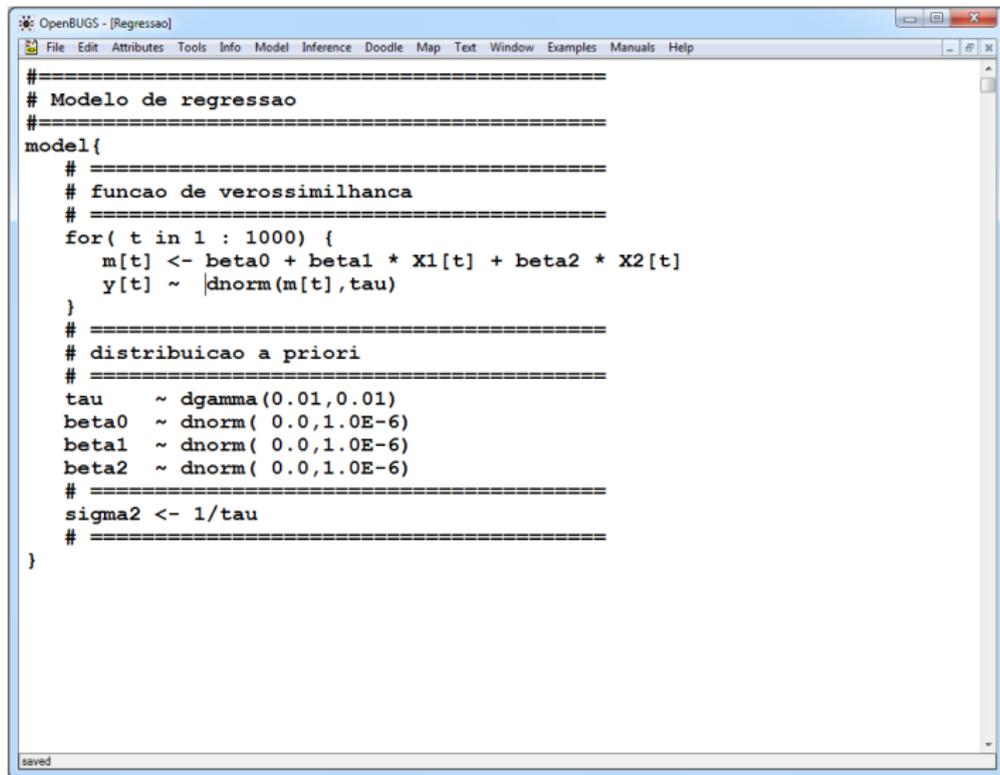
Dados artificiais gerados do seguinte modelo:

$$y_t = \beta_0 + \beta_1 x_{1t} + \beta_2 x_{2t} + \varepsilon_t, \quad t = 1, 2, \dots, n,$$

com $\beta_0 = 1, 0, \beta_1 = 5, 0, \beta_2 = -8, 0, \sigma^2 = 0, 25$ e $n = 1000$.

O OpenBUGS trabalha com a precisão $\phi = 1/\sigma^2$.

O modelo em código BUGS



```
OpenBUGS - [Regressao]
File Edit Attributes Tools Info Model Inference Doodle Map Text Window Examples Manuals Help

#####
# Modelo de regressao
#####
model{
  # =====
  # funcao de verossimilhanca
  # =====
  for( t in 1 : 1000) {
    m[t] <- beta0 + beta1 * X1[t] + beta2 * X2[t]
    y[t] ~ dnorm(m[t],tau)
  }
  # =====
  # distribuicao a priori
  # =====
  tau ~ dgamma(0.01,0.01)
  beta0 ~ dnorm( 0.0,1.0E-6)
  beta1 ~ dnorm( 0.0,1.0E-6)
  beta2 ~ dnorm( 0.0,1.0E-6)
  # =====
  sigma2 <- 1/tau
  # =====
}

saved
```

Opções básicas do OpenBUGS

1. **Model** → **Specification...**
2. **check model**
Verificar por *model is syntactically correct*
3. **load data**
Verificar por *data loaded*
4. **compile**
Verificar por *model compiled*
5. **load inits** ou **gen inits**
Verificar por *initial values generated, model initialized*
6. **Inference** → **Samples...**
7. **Model** → **Update...**

OpenBUGS

File Edit Attributes Tools Info Model Inference Doodle Map Text Window Examples Manuals Help

Regressao

```

=====
# Modelo de regressao
=====
model{
#
# funcao de verossimilhanca
=====
for( t in 1 : 1000) {
  m[t] <- beta0 + betal * X1[t] + beta2 * X2[t]
  y[t] ~ dnorm(m[t],tau)
}
#
# distribuicao a priori
=====
tau ~ dgamma( 0.01, 0.01)
beta0 ~ dnorm( 0.0, 1.0E-6)
betal ~ dnorm( 0.0, 1.0E-6)
beta2 ~ dnorm( 0.0, 1.0E-6)
#
sigma2 <- 1/tau
=====
}

```

dados_regressao.txt

```

y[] X1[] X2[]
1.5813 -1.0427 -0.6925
-9.2750 -0.1839 1.1247
0.2973 0.1793 0.1745
13.3930 -0.6686 -1.9678
9.6827 -1.0044 -1.6443
3.0409 0.1368 -0.1098
9.9864 0.8336 -0.6367
1.4518 -0.9519 -0.6803
4.4801 1.2304 0.3622
-4.4061 0.5661 0.9546
-0.3773 -0.3434 -0.0200
-4.5558 -0.4316 0.3574
6.6601 -0.2258 -0.8589
-2.1320 -0.5097 0.1315
8.3574 -0.5087 -1.2755
14.6847 0.7705 -1.2487
-2.2954 -1.1900 -0.3454
-5.1337 -0.9577 0.1353
0.8857 -1.3660 -0.7683
-15.3911 -2.7752 0.2730
4.9115 -2.1968 -1.7757
2.7169 0.4462 0.0154
-10.9538 -0.9457 0.9024
-5.5844 0.4429 1.0325
17.9904 1.8273 -0.8804
11.3244 1.0142 -0.6734
-7.8382 0.0416 1.0017

```

Specification Tool

check model load data

compile num of chans 2

load into for chan 1

gen into

Sample Monitor Tool

node - chans 1 to 2 percentiles

beg 10001 end 10000000 thin 5

clear set diagnostics trace jump

stats density bpr diag history accept

code quantiles auto cor

percentiles: 5 15 25 75 90 95 97.5

Update Tool

updates 50000 refresh 1000

update thin 1 iteration 50000

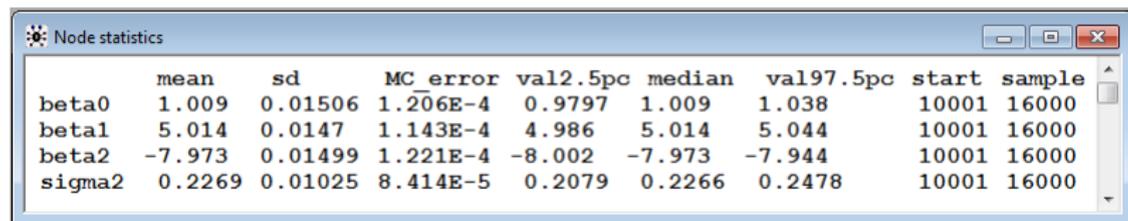
adapting over relax

Node statistics

	mean	sd	MC error	val0.5po	median	val97.5po	start	sample
beta0	1.009	0.01506	1.206E-4	0.9797	1.009	1.038	10001	16000
beta1	5.014	0.0147	1.143E-4	4.986	5.014	5.044	10001	16000
beta2	-7.973	0.01499	1.221E-4	-8.002	-7.973	-7.944	10001	16000
sigma2	0.2269	0.01025	8.414E-5	0.2079	0.2266	0.2478	10001	16000

saved

Resultados



	mean	sd	MC_error	val2.5pc	median	val97.5pc	start	sample
beta0	1.009	0.01506	1.206E-4	0.9797	1.009	1.038	10001	16000
beta1	5.014	0.0147	1.143E-4	4.986	5.014	5.044	10001	16000
beta2	-7.973	0.01499	1.221E-4	-8.002	-7.973	-7.944	10001	16000
sigma2	0.2269	0.01025	8.414E-5	0.2079	0.2266	0.2478	10001	16000

Parâmetro	Média	D.P.	2,5%	Mediana	97,5%
β_0	1,009	0,0151	0,980	1,009	1,039
β_1	5,014	0,0147	4,986	5,014	5,043
β_2	-7,973	0,0150	-8,002	-7,973	-7,944
σ^2	0,227	0,0103	0,208	0,227	0,248

Resultados

