Estatística e Análise de Dados em Zootecnia 2025/2026

Modelo Linear

Elsa Gonçalves Secção de Matemática, DCEB, ISA-Ulisboa

(Adaptado, Cadima J. (2021). O Modelo Linear, ISA, UILisboa)

Regressão Linear – Abordagem Inferencial

Regressão Linear - Inferência

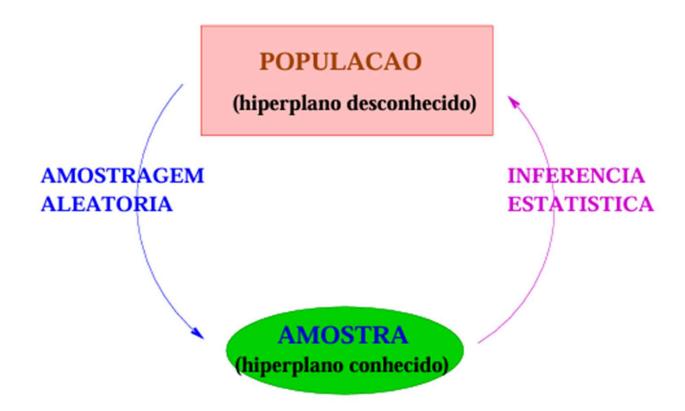
- Até aqui a regressão linear foi usada apenas como técnica descritiva. Se as n observações forem a totalidade da população de interesse, pouco mais há a dizer.
- Mas, com frequência, as n observações são apenas uma amostra aleatória de uma população maior.
- Um hiperplano ajustado a partir duma dada amostra,
 y = b₀ + b₁ x₁ + b₂ x₂ + ... + b_p x_p, é apenas uma estimativa de um hiperplano populacional

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_p x_p$$
.

Outras amostras dariam hiperplanos ajustados diferentes.

Coloca-se o problema da inferência estatística.

O problema da Inferência Estatística na Reg. Linear



MODELO - Regressão Linear

A fim de se poder fazer inferência sobre o hiperplano populacional, vamos admitir pressupostos adicionais.

- Y variável resposta aleatória.
- $x_1, ..., x_p$ variáveis preditoras não aleatórias (fixadas pelo experimentador ou trabalha-se condicionalmente aos valores de $x_1, ..., x_p$)

O modelo será ajustado com base em:

 $\{(x_{1(i)}, x_{2(i)}, ..., x_{p(i)}, Y_i)\}_{i=1}^n - n \text{ conjuntos de observações independentes das variáveis } x_1, x_2, ..., x_p e Y, \text{ sobre } n \text{ unidades experimentais.}$

MODELO RL – Linearidade

Vamos ainda admitir que a relação de fundo entre Y e x_1 , x_2 , ..., x_p , é linear (afim), com uma variabilidade aleatória em torno dessa relação, representada por um erro aleatório ε . Para todo o i = 1, ..., n:

$$Y_i = \beta_0 + \beta_1 \quad x_{1(i)} + ... + \beta_p \quad x_{p(i)} + \epsilon_i$$

 $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$

v.a. $cte.$ $cte.$ $cte.$ $cte.$ $cte.$ $cte.$ $cte.$ $cte.$ $cte.$

MODELO Regressão Linear – Os erros aleatórios

Vamos ainda admitir que os erros aleatórios ε_i :

Têm valor esperado (valor médio) nulo:

$$E[\varepsilon_i] = 0$$
, $\forall i = 1,...,n$

(não é hipótese restritiva).

- Têm distribuição Normal (é restritiva, mas bastante geral).
- Homogeneidade de variâncias: têm sempre a mesma variância

$$V[\varepsilon_i] = \sigma^2$$
, $\forall i = 1,...,n$

(é restritiva, mas conveniente).

 São variáveis aleatórias independentes (é restritiva, mas conveniente).

O Modelo Linear

O modelo para inferência na regressão linear é assim:

O Modelo Linear

- $Y_i = \beta_0 + \beta_1 x_{1(i)} + \beta_2 x_{2(i)} + \cdots + \beta_p x_{p(i)} + \varepsilon_i, \quad \forall i = 1, ..., n.$
- $\mathfrak{e}_i \sim \mathcal{N}(0, \sigma^2), \quad \forall i = 1, ..., n.$
- **3** $\{\varepsilon_i\}_{i=1}^n$ v.a. independentes.

NOTA: Os erros aleatórios são variáveis aleatórias independentes e identicamente distribuídas (i.i.d.).

Dado o modelo, o valor esperado (médio) de Y_i , condicional aos valores $x_1, x_2, ..., x_p$ dos preditores, é:

$$\mu_i = E[Y_i | x_1, x_2, ..., x_p] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_p x_p$$

NOTA: β_j ($j \neq 0$) é a variação média em Y, associada a um aumento de uma unidade em x_j , mantendo os restantes preditores constantes.

O estudo do modelo

Um primeiro objectivo da inferência: os p+1 parâmetros do modelo, β_j (j=0,1,...,p).

Os parâmetros ajustados $\vec{\mathbf{b}} = (b_0, b_1, b_2, ..., b_p)$, são estimativas desses parâmetros.

Para ser possível construir intervalos de confiança e/ou efectuar testes de hipóteses sobre os valores dos parâmetros populacionais β_i , há-que:

- Definir estimadores $\hat{\beta}_i$ dos parâmetros populacionais;
- conhecer as respectivas distribuições de probabilidades (ao abrigo do Modelo);

A validade da inferência depende da validade dos pressupostos do modelo.

A notação matricial/vectorial

$$\begin{array}{rcl} Y_1 & = & \beta_0 + \beta_1 x_{1(1)} + \beta_2 x_{2(1)} + \cdots + \beta_p x_{p(1)} & + & \varepsilon_1 \\ Y_2 & = & \beta_0 + \beta_1 x_{1(2)} + \beta_2 x_{2(2)} + \cdots + \beta_p x_{p(2)} & + & \varepsilon_2 \\ Y_3 & = & \beta_0 + \beta_1 x_{1(3)} + \beta_2 x_{2(3)} + \cdots + \beta_p x_{p(3)} & + & \varepsilon_3 \\ \vdots & \vdots & & \vdots & \vdots & \vdots \\ Y_n & = & \underbrace{\beta_0 + \beta_1 x_{1(n)} + \beta_2 x_{2(n)} + \cdots + \beta_p x_{p(n)}}_{=\vec{\mathbf{X}}\vec{\mathbf{B}}} & + & \underbrace{\varepsilon_n}_{=\vec{\mathbf{E}}} \end{array}$$

As n equações correspondem a uma única equação vectorial:

$$\vec{Y} = X\vec{\beta} + \vec{\epsilon} ,$$

onde:

$$\vec{\mathbf{Y}} = \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ Y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_{1_{(1)}} & x_{2_{(1)}} & \cdots & x_{\rho_{(1)}} \\ 1 & x_{1_{(2)}} & x_{2_{(2)}} & \cdots & x_{\rho_{(2)}} \\ 1 & x_{1_{(3)}} & x_{2_{(3)}} & \cdots & x_{\rho_{(3)}} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1_{(n)}} & x_{2_{(n)}} & \cdots & x_{\rho_{(n)}} \end{bmatrix}, \quad \vec{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}, \quad \vec{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

- \vec{Y} e $\vec{\epsilon}$ são vectores aleatórios,
- X é uma matriz não aleatória e β um vector não-aleatório.

Modelo Regressão Linear - versão vectorial

O Modelo Linear em notação vectorial

- $\vec{\epsilon} \sim \mathcal{N}_{n}(\vec{0}, \sigma^{2} \mathbf{I}_{n}), \text{ com } \vec{0} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \dots \\ 0 \end{bmatrix}; \sigma^{2} \mathbf{I}_{n} = \begin{bmatrix} 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & \sigma^{2} & 0 & \dots & 0 \\ 0 & 0 & \sigma^{2} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma^{2} \end{bmatrix}$
- Cada erro aleatório individual ε_i tem distribuição Normal.
- Cada erro aleatório individual tem média zero: $E[\varepsilon_i] = 0$.
- Cada erro aleatório individual tem variância igual: $V[\varepsilon_i] = \sigma^2$.
- Erros aleatórios diferentes são independentes, porque Cov[ε_i, ε_j] = 0 se i ≠ j e, numa Multinormal, isso implica a independência.

A distribuição de **Y**

Teorema (Primeiras Consequências do Modelo)

Dado o Modelo de Regressão Linear, tem-se:

$$\vec{\mathbf{Y}} \sim \mathcal{N}_n(\mathbf{X}\vec{\boldsymbol{\beta}}, \sigma^2 \mathbf{I}_n)$$

De facto, \vec{Y} é soma de vector não aleatório $(\vec{x}\vec{\beta})$ e vector aleatório $(\vec{\epsilon})$:

$$\vec{\mathbf{Y}} = \underbrace{\mathbf{X}\vec{\boldsymbol{\beta}}}_{="\vec{\mathbf{Z}}"} + \underbrace{\vec{\boldsymbol{\varepsilon}}}_{="\vec{\mathbf{Z}}"}.$$

- $\vec{\epsilon} \sim \mathcal{N}(\vec{\mathbf{0}}, \sigma^2 \mathbf{I}_n)$.
- Somar vector constante $(\mathbf{X}\vec{\boldsymbol{\beta}})$ a um vector aleatório Multinormal $(\vec{\boldsymbol{\varepsilon}})$ não destrói a Multinormalidade.
- $E[\vec{Y}] = E[X\vec{\beta} + \vec{\epsilon}] = X\vec{\beta} + E[\vec{\epsilon}] = X\vec{\beta}$.
- $V[\vec{\mathbf{Y}}] = V[\mathbf{X}\vec{\boldsymbol{\beta}} + \vec{\boldsymbol{\varepsilon}}] = V[\vec{\boldsymbol{\varepsilon}}] = \sigma^2 \mathbf{I}_n$

A distribuição de **Y** (interpretação)

$$\vec{\mathbf{Y}} \sim \mathscr{N}_n(\mathbf{X}\vec{\boldsymbol{\beta}}, \sigma^2 \mathbf{I}_n)$$

Tendo em conta as propriedades da Multinormal:

- Cada observação individual Y_i tem distribuição Normal.
- Cada observação individual Y_i tem média $\mu_i = E[Y_i] = \vec{\mathbf{x}}_{[i,]}^t \vec{\boldsymbol{\beta}} = \beta_0 + \beta_1 x_{1(i)} + \beta_2 x_{2(i)} + ... + \beta_p x_{p(i)}.$
- Cada observação individual tem variância igual: $V[Y_i] = \sigma^2$.
- Observações diferentes de Y são independentes, porque Cov[Y_i, Y_j] = 0 se i ≠ j e, numa Multinormal, isso implica a independência.

O vector de estimadores $\hat{\beta}$

O vector de estimadores $\vec{\hat{\beta}} = (\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p)^t$ é definido a partir da equação do vector \vec{b} de estimativas mas substituindo o vector \vec{y} de valores observados de Y pelo vector aleatório \vec{Y} .

Estimadores de Mínimos Quadrados dos parâmetros

$$\hat{\hat{\beta}} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \vec{\mathbf{Y}}.$$

O vector $\vec{\hat{\beta}}$ é de dimensão p+1. O seu primeiro elemento é o estimador de β_0 , o seu segundo elemento é o estimador de β_1 , etc... Em geral, o estimador de β_j está na posição j+1 do vector $\vec{\hat{\beta}}$.

Os resultados gerais já referidos permitem facilmente determinar a distribuição de probabilidades do estimador $\hat{\beta}$.

A distribuição do vector de estimadores $\hat{oldsymbol{eta}}$

Teorema (Distribuição do estimador $\hat{\beta}$)

Dado o Modelo de Regressão Linear Múltipla, tem-se:

$$\vec{\hat{\beta}} \sim \mathscr{N}_{p+1}(\vec{\beta}, \sigma^2(\mathbf{X}^t\mathbf{X})^{-1})$$

 $\hat{\boldsymbol{\beta}}$ é produto de matriz não aleatória, $(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t$, e vector aleatório, $\vec{\mathbf{Y}}$:

$$\hat{\boldsymbol{\beta}} = \underbrace{(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t}_{\mathbf{B}''} \underbrace{\vec{\mathbf{Y}}}_{\mathbf{Z}''}.$$

- $\vec{\mathbf{Y}} \sim \mathcal{N}_n(\mathbf{X}\vec{\boldsymbol{\beta}}, \sigma^2\mathbf{I}_n)$.
- Multiplicar matriz constante, (X^tX)⁻¹X^t, por um vector aleatório Multinormal (Y)
 não destrói a Multinormalidade.
- $\mathbf{E}[\vec{\hat{\beta}}] = \mathbf{E}[(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \vec{\mathbf{Y}}] = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{E}[\vec{\mathbf{Y}}] = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{X} \vec{\beta} = \vec{\beta}.$
- $V[\hat{\boldsymbol{\beta}}] = V[(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \hat{\mathbf{Y}}] = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t V[\hat{\mathbf{Y}}][(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t]^t = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \cdot \sigma^2 \mathbf{I}_n \cdot \mathbf{X}[(\mathbf{X}^t \mathbf{X})^{-1}]^t = \sigma^2 \cdot (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{X}[(\mathbf{X}^t \mathbf{X})^t]^{-1} = \sigma^2 (\mathbf{X}^t \mathbf{X})^{-1}.$

A distribuição de $\vec{\hat{\beta}}$ (interpretação)

$$\vec{\hat{\beta}} \sim \mathscr{N}_{p+1}(\vec{\beta}, \sigma^2(\mathbf{X}^t\mathbf{X})^{-1})$$
.

Tendo em conta as propriedades da Multinormal

- Cada estimador individual β̂_i tem distribuição Normal.
- Cada estimador individual tem média $E[\hat{\beta}_j] = \beta_j$, logo é centrado
- Cada estimador individual tem variância $V[\hat{\beta}_j] = \sigma^2 (\mathbf{X}^t \mathbf{X})^{-1}_{(j+1,j+1)}$. (Note-se o desfasamento nos índices).
- Estimadores individuais diferentes não são (em geral) independentes, porque $(\mathbf{X}^t\mathbf{X})^{-1}$ não é, em geral, uma matriz diagonal: $Cov[\hat{\beta}_i, \hat{\beta}_j] = \sigma^2 (\mathbf{X}^t\mathbf{X})^{-1}_{(i+1,i+1)}$.
- Logo, o estimador $\hat{\beta}_j$ de um parâmetro individual β_j tem distribuição $\hat{\beta}_j \sim \mathcal{N}(\beta_j, \sigma_{\hat{\beta}_i}^2)$, com $\sigma_{\hat{\beta}_i}^2 = \sigma^2(\mathbf{X}^t\mathbf{X})_{(j+1,j+1)}^{-1}$.

Estimação dos parâmetros do Modelo RLS

A recta do modelo RLS tem dois parâmetros: β_0 e β_1 .

Definem-se estimadores desses parâmetros a partir das expressões amostrais obtidas para b_0 e b_1 pelo Método dos Mínimos Quadrados.

Recordar:
$$b_1 = \frac{cov_{xy}}{s_x^2} = \frac{\sum\limits_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})}{(n-1) s_x^2} = \frac{\sum\limits_{i=1}^{n} (x_i - \overline{x})y_i}{(n-1) s_x^2} = \sum\limits_{i=1}^{n} \frac{x_i - \overline{x}}{(n-1) s_x^2} y_i$$

Estimador de β_1

$$\hat{\beta}_1 = \sum_{i=1}^n \frac{x_i - \overline{x}}{\frac{(n-1)}{S_X^2}} Y_i = \sum_{i=1}^n c_i Y_i, \quad \text{com } c_i = \frac{x_i - \overline{x}}{\frac{(n-1)}{S_X^2}}$$

Nota: O estimador $\hat{\beta}_1$ é combinação linear de Normais independentes, logo tem distribuição Normal.

Estimação dos parâmetros do Modelo RLS (cont.)

Recordar: $b_0 = \overline{y} - b_1 \overline{x}$.

Estimador de β_0

$$\hat{\beta}_0 = \overline{Y} - \hat{\beta}_1 \overline{x} = \frac{1}{n} \sum_{i=1}^n Y_i - \overline{x} \sum_{i=1}^n c_i Y_i = \sum_{i=1}^n \left(\frac{1}{n} - \overline{x} c_i \right) Y_i = \sum_{i=1}^n d_i Y_i,$$

com

$$d_i = \frac{1}{n} - \overline{x}c_i = \frac{1}{n} - \frac{(x_i - \overline{x})\overline{x}}{(n-1)S_x^2}.$$

Quer $\hat{\beta}_1$, quer $\hat{\beta}_0$, são combinações lineares das observações $\{Y_i\}_{i=1}^n$, logo são combinações lineares de variáveis aleatórias Normais independentes. Logo, ambos os estimadores têm distribuição Normal.

Distribuição dos estimadores RLS

Distribuição dos estimadores dos parâmetros

Dado o Modelo de Regressão Linear Simples,

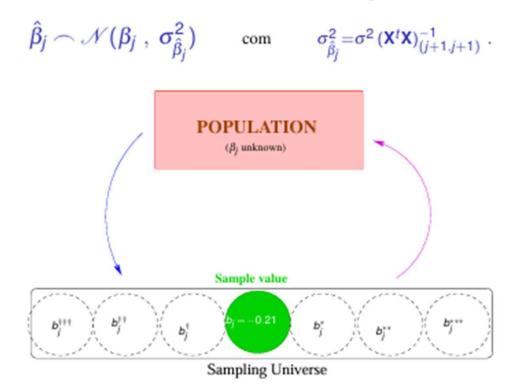
$$\hat{\beta}_1 \sim \mathcal{N}\left(\beta_1, \frac{\sigma^2}{(n-1)S_x^2}\right),$$

$$\hat{\beta}_0 \sim \mathcal{N}\left(\beta_0, \sigma^2 \left[\frac{1}{n} + \frac{\overline{x}^2}{(n-1)S_x^2}\right]\right)$$

NOTAS:

- **1** Ambos os estimadores são centrados: $E[\hat{\beta}_1] = \beta_1$ e $E[\hat{\beta}_0] = \beta_0$.
- Quanto maior (n-1) S_X^2 , menor a variância dos estimadores.
- **3** A variância de $\hat{\beta_0}$ também diminui com o aumento de n, e com a maior proximidade de \overline{x} de zero.

A distribuição na amostragem de $\hat{\beta}_i$ (interpretação)



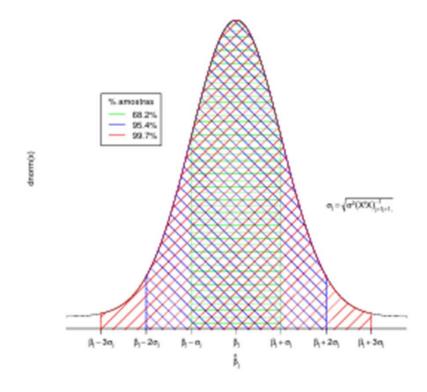
O conjunto de todas as possíveis amostras de dimensão *n* designa-se o Universo de Amostragem

A distribuição de probabilidades de $\hat{\beta}_j$ pode ser vista como a distribuição dos valores de b_i ao longo do Universo de Amostragem.

A distribuição na amostragem de $\hat{\beta}_j$ (interpretação)

$$\hat{\beta}_j \frown \mathcal{N}(\beta_j, \sigma_{\hat{\beta}_j}^2)$$
 com $\sigma_{\hat{\beta}_j}^2 = \sigma^2(\mathbf{X}^t\mathbf{X})_{(j+1,j+1)}^{-1}$.

Distribuição na amostragem de $\hat{\beta}$



A distribuição dum estimador individual

Como se viu, tem-se, $\forall j = 0, 1, ..., p$:

$$\hat{\beta}_{j} \quad \frown \quad \mathscr{N}\left(\beta_{j} , \sigma^{2}(\mathbf{x}^{t}\mathbf{x})_{(j+1,j+1)}^{-1}\right)$$

$$\Leftrightarrow \quad \frac{\hat{\beta}_{j} - \beta_{j}}{\sigma_{\hat{\beta}_{j}}} \quad \frown \quad \mathscr{N}(0,1) ,$$

$$\text{com } \sigma_{\hat{\beta}_j} = \sqrt{\sigma^2 (\mathbf{X}^t \mathbf{X})_{(j+1,j+1)}^{-1}}.$$

Este resultado distribucional permitiria construir intervalos de confiança ou fazer testes a hipóteses sobre os parâmetros $\vec{\beta}$, não fosse o desconhecimento da variância σ^2 dos erros aleatórios.

Distribuição dos estimadores RLS

Distribuição dos estimadores (cont.)

Dado o Modelo de Regressão Linear Simples,

NOTAS:

- O desvio padrão dum estimador designa-se erro padrão (em inglês, standard error).
- Não confundir os erros padrão dos estimadores, $\sigma_{\hat{\beta}_1}$ e $\sigma_{\hat{\beta}_0}$, com o desvio padrão σ dos erros aleatórios.

O problema de σ^2 desconhecido

Para poder utilizar um estimador $\hat{\beta}_j$ na inferência, é preciso conhecer a sua distribuição de probabilidades, sem a presença de quantidades não-amostrais desconhecidas, além de β_j .

Para ultrapassar este problema é preciso:

- obter um estimador para σ^2 ; e
- ver o que acontece à distribuição de $\hat{\beta}_j$ quando σ^2 é substituído pelo seu estimador.

Como $\sigma^2 = V(\varepsilon_i)$, $\forall i$, e como os erros aleatórios ε_i são desconhecidos, é natural procurar um estimador de σ^2 através dos resíduos.

Estimando σ^2

Erros aleatórios (variáveis aleatórias – não observáveis)
$$\varepsilon_i = Y_i - (\beta_0 + \beta_1 x_{1(i)} + \beta_2 x_{2(i)} + ... + \beta_p x_{p(i)})$$
 Resíduos (variáveis aleatórias – observáveis)
$$E_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{1(i)} + \hat{\beta}_2 x_{2(i)} + ... + \hat{\beta}_p x_{p(i)})$$

$$= \hat{Y}_i$$
 Resíduos (observados)
$$e_i = y_i - (b_0 + b_1 x_{1(i)} + b_2 x_{2(i)} + ... + b_p x_{p(i)})$$

Quadrado Médio Residual (QMRE)

Define-se o Quadrado Médio Residual como

QMRE =
$$\frac{SQRE}{n-(p+1)} = \frac{\sum_{i=1}^{n} E_i^2}{n-(p+1)}$$

Dado o Modelo Linear, $\hat{\sigma}^2 = QMRE$ é um estimador centrado da variância comum dos erros aleatórios, $\sigma^2 = V[\varepsilon_i]$:

$$E[QMRE] = \sigma^2$$
.

O Quadrado Médio Residual tem como unidades de medida o quadrado das unidades de Y.

Quantidades fulcrais para a inferência sobre β_j

A estimação dos erros padrão com o QMRE transforma as distribuições normais em distribuições *t-Student*

Teorema (Distribuições para a inferência sobre β_i)

Dado o Modelo de Regressão Linear Múltipla, tem-se

$$\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma}_{\hat{\beta}_j}} \frown t_{n-(p+1)}, \qquad \forall j = 0, 1, ..., p$$

$$com \ \hat{\sigma}_{\hat{\beta}_j} = \sqrt{QMRE \cdot (\mathbf{X}^t \mathbf{X})_{(j+1,j+1)}^{-1}}.$$

Este Teorema dá-nos os resultados que servem de base à construção de intervalos de confiança e testes de hipóteses para os parâmetros β_j do modelo populacional.

Quantidades centrais para a inferência sobre β_0 e β_1

A estimação dos erros padrão com o QMRE transforma as distribuições normais em distribuições *t-Student*

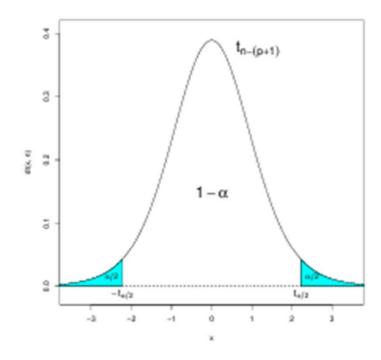
Distribuições *t-Student* para a inferência sobre β_0 e β_1

Dado o Modelo de Regressão Linear Simples, prova-se que:

Este Teorema é crucial, pois dá-nos os resultados que servirão de base à construção de intervalos de confiança e testes de hipóteses para os parâmetros da recta populacional, β_0 e β_1 .

Dedução de intervalo de confiança para β_i

Sabemos que $\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma}_{\hat{\beta}_j}} \ \frown \ t_{n-(p+1)}$. Logo,



$$P\left[-t_{\frac{\alpha}{2}} < \frac{\hat{\beta}_{j} - \beta_{j}}{\hat{\sigma}_{\hat{\beta}_{j}}} < t_{\frac{\alpha}{2}}\right] = 1 - \alpha$$

Dedução IC para β_i (cont.)

Trabalhar a dupla desigualdade até isolar β_i :

$$P\left[-t_{\frac{\alpha}{2}} < \frac{\hat{\beta}_{j}-\beta_{j}}{\hat{\sigma}_{\hat{\beta}_{j}}} < t_{\frac{\alpha}{2}}\right] = 1-\alpha$$

$$\begin{aligned} -t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\beta}_{j}} &< \hat{\beta}_{j} - \beta_{j} &< t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\beta}_{j}} \\ \Leftrightarrow & t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\beta}_{j}} > \beta_{j} - \hat{\beta}_{j} &> -t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\beta}_{j}} \\ \Leftrightarrow & \hat{\beta}_{j} - t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\beta}_{j}} &< \beta_{j} &< \hat{\beta}_{j} + t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\beta}_{j}} \end{aligned}$$

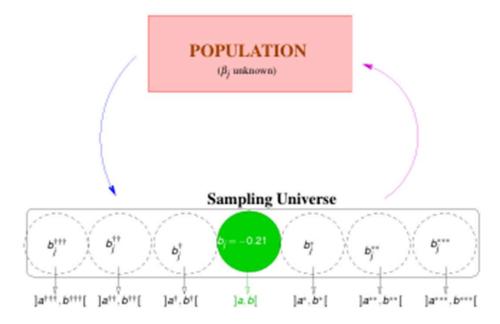
O intervalo aleatório

$$] \hat{\beta}_j - t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\beta}_j} , \hat{\beta}_j + t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\beta}_j} [$$

contém β_i com probabilidade $1 - \alpha$.

Intervalo aleatório para β_i (interpretação)

$$] \hat{\beta}_j - t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\beta}_j} , \hat{\beta}_j + t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\beta}_j} [$$



Cada amostra no Universo de Amostragem gera um intervalo concreto, chamado Intervalo de Confiança

Uma proporção $1-\alpha$ desses intervalos contêm o verdadeiro valor de β_j . Os restantes α não contêm β_j .

Intervalo de confiança para β_i

Intervalo de Confiança a $(1 - \alpha) \times 100\%$ para β_i

Dado o Modelo de Regressão Linear Múltipla e uma amostra, eis o intervalo a $(1-\alpha) \times 100\%$ de confiança para o parâmetro β_i :

sendo:

- b_j o elemento j+1 do vector das estimativas $\vec{\mathbf{b}}$
- $t_{\frac{\alpha}{2}[n-(p+1)]}$ o quantil de ordem $1-\frac{\alpha}{2}$ da distribuição $t_{n-(p+1)}$;
- $\hat{\sigma}_{\hat{\beta}_j} = \sqrt{QMRE \cdot (\mathbf{X}^t \mathbf{X})_{(j+1,j+1)}^{-1}}$ (com o valor de QMRE na nossa amostra).

NOTA: A amplitude do IC aumenta com *QMRE* e o valor diagonal da matriz $(\mathbf{X}^t\mathbf{X})^{-1}$ correspondente ao parâmetro β_i .

Intervalo de confiança para β_1

Intervalo de Confiança a $(1-\alpha) \times 100\%$ para β_1

Dado o Modelo RLS, um intervalo a $(1-\alpha) \times 100\%$ de confiança para o declive β_1 da recta de regressão populacional é dado por:

$$\ \, \Big] \, b_1 - t_{\frac{\alpha}{2}[n-2]} \, \hat{\sigma}_{\hat{\beta_1}} \quad , \quad b_1 + t_{\frac{\alpha}{2}[n-2]} \, \hat{\sigma}_{\hat{\beta_1}} \, \Big[\ ,$$

tendo $t_{\frac{\alpha}{2}[n-2]}$, b_1 e $\hat{\sigma}_{\hat{\beta_1}}$ sido definidos em acetatos anteriores.

NOTAS:

- O intervalo é centrado em b₁.
- A amplitude do intervalo é $2 \times t_{\frac{\alpha}{2}[n-2]} \hat{\sigma}_{\hat{\beta_1}}$.
- A amplitude aumenta com *QMRE* e diminui com $n \in S_X^2$: $\hat{\sigma}_{\hat{\beta}_1} = \sqrt{\frac{QMRE}{(n-1)S_X^2}}$
- A amplitude do IC aumenta para maiores graus de confiança $1-\alpha$.

Intervalo de confiança para β_0

Intervalo de Confiança a $(1-\alpha) \times 100\%$ para β_0

Dado o Modelo de Regressão Linear Simples, um intervalo a $(1-\alpha) \times 100\%$ de confiança para a ordenada na origem, β_0 , da recta populacional é:

$$\left] b_0 - t_{\frac{\alpha}{2}[n-2]} \cdot \hat{\sigma}_{\hat{\beta_0}} \quad , \quad b_0 + t_{\frac{\alpha}{2}[n-2]} \cdot \hat{\sigma}_{\hat{\beta_0}} \right[,$$

onde $t_{\frac{\alpha}{2}[n-2]}$, b_0 e $\hat{\sigma}_{\hat{\beta}_0}$ foram definidos em acetatos anteriores.

NOTAS:

- O intervalo é centrado em b₀.
- A amplitude do intervalo é $2 \times t_{\frac{n}{2}[n-2]} \hat{\sigma}_{\hat{\beta_0}}$.
- A amplitude aumenta com QMRE e com \overline{x}^2 e diminui com $n \in s_x^2$:

$$\hat{\sigma}_{\hat{\beta}_0} = \sqrt{QMRE \cdot \left[\frac{1}{n} + \frac{\overline{x}^2}{\frac{(n-1)}{S_X^2}}\right]}$$

A amplitude do IC aumenta para maiores graus de confiança 1-α.

Ainda o exemplo dos lírios

RLM

□ proc reg data=iris;

model PetalWidth = SepalLength SepalWidth PetalLength/clb;

run;

Parameter Estimates							
Variable	DF	Parameter Estimate	And the second second second second second	t Value	Pr > t	95% Confidence Limits	
Intercept	1	-0.24031	0.17837	-1.35	0.1800	-0.59283	0.11221
SepalLength	1	-0.20727	0.04751	-4.36	<.0001	-0.30115	-0.11338
SepalWidth	1	0.22283	0.04894	4.55	<.0001	0.12611	0.31955
PetalLength	1	0.52408	0.02449	21.40	<.0001	0.47568	0.57249

Exemplo b_1 : na nossa amostra estima-se que, em média, a largura da pétala diminui 0.20727 cm por cada aumento de 1 cm no comprimentos da sépala (mantendo-se as outras medições constantes).

Como $t_{0.025(146)=1.976346}$, o IC a 95% para β_1 é:

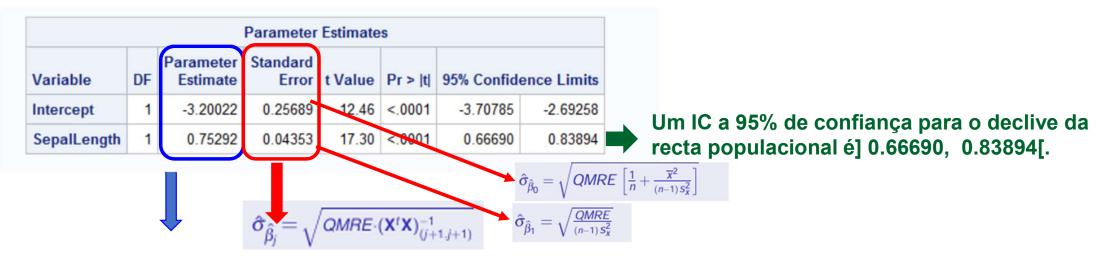
 $](-0.20727) - (1.976346)(0.04751) \ , (-0.20727) + (1.976346)(0.04751)[$

 \Leftrightarrow]-0.30115, -0.11338[

As estimativas dos desvios padrão associados à estimação de cada um dos parâmetros

Temos 95% de confiança em como o verdadeiro valor de β_1 (na população) está compreendido entre -0.30115 e -0.11338.

Ainda o exemplo dos lírios



Nota: O coeficiente associado ao preditor *Sepal.Length* na regressão linear simples agora ajustada é positivo, b_1 =0.75292. No modelo de regressão linear múltipla obteve-se um resultado differente, pois contém, além do preditor comprimento da sépala, outros dois preditores (largura da sépala e comprimento da pétala), que contribuem para a formação dos valores ajustados. Na presença desses dois preditores, a contribuição do comprimento da sépala teve um sinal negativo. Esta aparente contradição sublinha uma ideia importante: a introdução (ou exclusão) de preditores numa regressão linear têm efeitos sobre todos os parâmetros, não sendo possível prever qual será a equação ajustada sem refazer as contas do ajustamento.

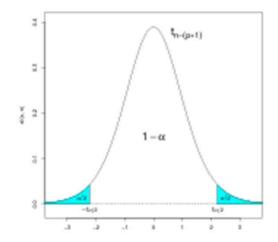
Testes de Hipóteses sobre os parâmetros

O resultado usado para construir ICs também permite Testes a Hipóteses sobre cada β_i . Admitindo a Hipótese Nula $H_0: \beta_i = c$:

$$T = \frac{\hat{\beta}_j - \overbrace{\beta_{j|H_0}}^{=c}}{\hat{\sigma}_{\hat{\beta}_j}} \quad f_{n-(p+1)}, \quad \forall j = 0, 1, ..., p$$

Rejeita-se H_0 em favor da Hipótese Alternativa $H_1: \beta_j \neq c$ se o valor calculado de T na amostra, T_{calc} , recair numa das caudas da distribuição.

Fixando o Nível de Significância α , tem-se a Região Crítica:



Testes de Hipóteses (bilateral) a $\hat{\beta}_j$

Testes de Hipóteses a β_i (Modelo de Regressão Linear Múltipla)

Hipóteses: H_0 : $\beta_j = c$ vs. H_1 : $\beta_j \neq c$

Estatística do Teste:
$$T = \frac{\hat{\beta}_j - \widehat{\beta_j}|_{H_0}}{\hat{\sigma}_{\hat{\beta}_j}} \frown t_{n-(p+1)}$$
, se H_0 verdade.

Nível de significância do teste: α

Região Crítica (Região de Rejeição bilateral): Rejeitar Ho se

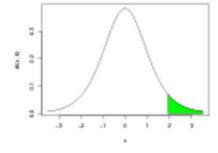
$$T_{calc} > t_{\frac{\alpha}{2}[n-(p+1)]}$$
 ou $T_{calc} < -t_{\frac{\alpha}{2}[n-(p+1)]}$

$$\iff$$
 $|T_{calc}| > t_{\frac{\alpha}{2}[n-(p+1)]}$

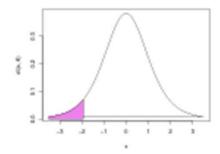
Testes de Hipóteses a $\hat{\beta}_i$ (unilaterais)

$$T = \frac{\hat{\beta}_j - \widehat{\beta_{j|H_0}}}{\hat{\sigma}_{\hat{\beta}_j}} - t_{n-(p+1)}$$

Com a Hipótese Alternativa $H_1: \beta_j > c$, só valores grandes da estatística sugerem a rejeição de $H_0: \beta_j \leq c$ (ou $H_0: \beta_j = c$):



Com a Hipótese Alternativa $H_1: \beta_j < c$, só valores pequenos de T_{calc} sugerem rejeitar $H_0: \beta_j \geq c$ (ou $H_0: \beta_j = c$):



Testes de Hipóteses sobre os parâmetros

Dado o Modelo de Regressão Linear Múltipla,

```
Testes de Hipóteses a \beta_i (Regressão Linear Múltipla)
 Hipóteses: H_0: \beta_j = c \text{ vs. } H_1: \beta_j \neq c
Estatística do Teste: T = \frac{\hat{\beta}_j - \widehat{\beta}_j|_{H_0}}{\hat{\sigma}_{\hat{\beta}_j}} \frown t_{n-(p+1)}, se H_0 verdade.
Nível de significância do teste: \alpha
Região Crítica (Região de Rejeição): Rejeitar H<sub>0</sub> se
                                                       (Unilateral esquerdo)
                   T_{calc} < -t_{\alpha[n-(p+1)]}
                    |T_{calc}| > t_{\alpha/2[n-(p+1)]}
                                                       (Bilateral)
                    T_{calc} > t_{\alpha[n-(p+1)]}
                                                       (Unilateral direito)
```

Os p-values

Valores de prova (*p-value*)

O *p-value* é a probabilidade da estatística de teste tomar valores mais extremos que o valor calculado a partir da amostra, sob H_0

O cálculo do *p-value* é feito de forma diferente, consoante a natureza da Região Crítica (RC) (unilateral direita ou esquerda, ou bilateral).

Sendo
$$T \sim t_{n-(p+1)}$$

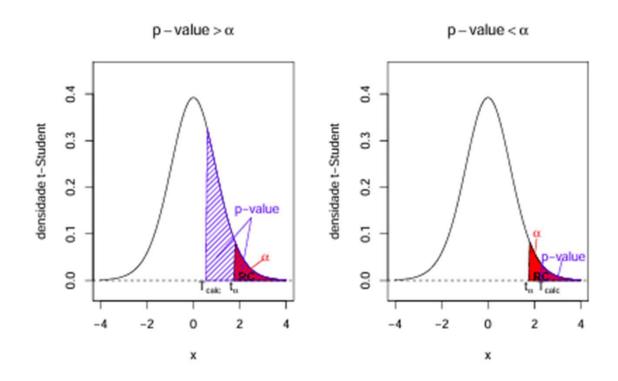
RC Unilateral direita: $p = P[T > T_{calc}]$

RC Unilateral esquerda: $p = P[T < T_{calc}]$

RC Bilateral: $p = 2 \times P[T > |T_{calc}|].$

A relação de *p-values* e níveis de significância

- p-value $> \alpha \Rightarrow$ não rejeição de H_0 ao nível α ;
- p−value < α ⇒ rejeição de H₀ ao nível α;



Em geral: p-value muito pequeno implica rejeição Ho.

RLM

proc reg data=iris;

model PetalWidth = SepalLength SepalWidth PetalLength/clb;

run;

			Parameter					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confider	nce Limits	
Intercept	1	-0.24031	0.17837	-1.35	0.1800	-0.59203	0.11221	→ Teste H_0 : $β_o$ =0 vs. H_1 : $β_o ≠0$
SepalLength	1	-0.20727	0.04751	-4.36	<.0001	0.30115	0.11338	Teste H_0 : β_1 =0 vs. H_1 : $\beta_1 \neq 0$
SepalWidth	1	0.22283	0.04894	4.55	<.0001	0.12611	0.31955	Teste H_0 : β_2 =0 vs. H_1 : $\beta_2 \neq 0$
PetalLength	1	0.52408	0.02449	21.40	<.0001	0.47568	0.57249	\rightarrow Teste H_0 : β_3 =0 vs. H_1 : $\beta_3 \neq 0$

Exemplo:
$$T_{Calc} = \frac{b_3 - \beta_3 | H_0}{\widehat{\sigma}_{\widehat{\beta}_3}} = \frac{0.52408}{0.02449} = 21.40$$

O valor de prova (*p-value*) indica uma claríssima rejeição da hipótese nula para um nível de significância usual

Nota: por exemplo, para o teste H_0 : β_3 =0.5 vs. H_1 : $\beta_3 \neq 0.5$

$$T_{Calc} = \frac{b_3 - \beta_3 | H_0}{\widehat{\sigma}_{\widehat{\beta}_3}} = \frac{0.52408 - 0.5}{0.02449} = 0.983258473$$

$$t_{0.05}_{146)} = t_{0.025(146)} \approx 1.96$$

$$|T_{Calc}| < 1.96, para \ \alpha = 0.05, n\~{a}o \ se \ rejeita \ a \ hip\'otese \ nula$$

(O valor de prova (p-value) da tabela não é válido neste caso)

Combinações lineares dos parâmetros

Seja $\vec{\mathbf{a}} = (a_0, a_1, ..., a_p)^t$ um vector não aleatório em \mathbb{R}^{p+1} . O produto interno $\vec{\mathbf{a}}^t \vec{\boldsymbol{\beta}}$ define uma combinação linear dos parâmetros do modelo:

$$\vec{a}^t \vec{\beta} = a_0 \beta_0 + a_1 \beta_1 + a_2 \beta_2 + ... + a_p \beta_p$$
.

Casos particulares importantes são se:

- \vec{a} tem um único elemento não-nulo, $a_{j+1} = 1$: $\vec{a}^t \vec{\beta} = \beta_j$.
- \vec{a} só tem dois elementos não-nulos, $a_{i+1} = 1$ e $a_{j+1} = \pm 1$: $\vec{a}^t \vec{\beta} = \beta_i \pm \beta_j$.
- $\vec{\mathbf{a}} = (1, x_1, x_2, ..., x_p)$: $\vec{\mathbf{a}}^t \vec{\boldsymbol{\beta}}$ é o valor esperado de Y associado aos valores indicados das variáveis preditoras:

$$\vec{\mathbf{a}}^{t}\vec{\boldsymbol{\beta}} = \beta_{0} + \beta_{1}x_{1} + \beta_{2}x_{2} + ... + \beta_{p}x_{p}$$

$$= E[Y | X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{p} = x_{p}]$$

$$= \mu_{Y | \vec{\mathbf{x}}}$$

Inferência sobre combinações lineares dos β_i s

Estima-se $\vec{a}^t \vec{\beta}$ com a mesma combinação linear dos estimadores:

$$\vec{a}^t \hat{\hat{\beta}} = a_0 \hat{\beta_0} + a_1 \hat{\beta_1} + a_2 \hat{\beta_2} + ... + a_p \hat{\beta_p}$$

Sabemos determinar a distribuição de probabilidades de $\vec{a}^t \hat{\beta}$:

- Sabemos que $\vec{\hat{\beta}} \sim \mathcal{N}_{p+1} \left(\vec{\beta}, \sigma^2 (\mathbf{X}^t \mathbf{X})^{-1} \right)$
- Logo, $\vec{\mathbf{a}}^t \hat{\vec{\boldsymbol{\beta}}} \sim \mathcal{N}_1(\vec{\mathbf{a}}^t \vec{\boldsymbol{\beta}}, \sigma^2 \vec{\mathbf{a}}^t (\mathbf{X}^t \mathbf{X})^{-1} \vec{\mathbf{a}})$
- Ou seja, $\vec{\mathbf{Z}} = \frac{\vec{\mathbf{a}}^t \vec{\hat{\boldsymbol{\beta}}} \vec{\mathbf{a}}^t \vec{\boldsymbol{\beta}}}{\sqrt{\sigma^2 \vec{\mathbf{a}}^t (\mathbf{X}^t \mathbf{X})^{-1} \vec{\mathbf{a}}}} \frown \mathcal{N}(0,1);$
- Por um raciocínio análogo ao usado nos β_j individuais, tem-se:

$$\frac{\vec{\mathbf{a}}^t \vec{\hat{\boldsymbol{\beta}}} - \vec{\mathbf{a}}^t \vec{\boldsymbol{\beta}}}{\sqrt{QMRE \cdot \vec{\mathbf{a}}^t (\mathbf{X}^t \mathbf{X})^{-1} \vec{\mathbf{a}}}} \frown t_{n-(p+1)}.$$

Quantidades centrais para a inferência sobre $\vec{a}^t \vec{\beta}$

Teorema (Distribuições para combinações lineares dos β s)

Dado o Modelo de Regressão Linear Múltipla, tem-se

$$\frac{\vec{\mathbf{a}}^t \hat{\vec{\boldsymbol{\beta}}} - \vec{\mathbf{a}}^t \vec{\boldsymbol{\beta}}}{\hat{\sigma}_{\vec{\mathbf{a}}^t \hat{\vec{\boldsymbol{\beta}}}}} \ \frown \ t_{n-(p+1)} \ ;$$

$$com \ \hat{\sigma}_{\vec{\mathbf{a}}^t \vec{\hat{\mathbf{B}}}} = \sqrt{QMRE \cdot \vec{\mathbf{a}}^t (\mathbf{X}^t \mathbf{X})^{-1} \vec{\mathbf{a}}}.$$

Este Teorema dá-nos os resultados que servem de base à construção de intervalos de confiança e testes de hipóteses para quaisquer combinações lineares dos parâmetros β_i do modelo.

Intervalo de confiança para $\vec{a}^t \vec{\beta}$

Intervalo de Confiança a $(1 - \alpha) \times 100\%$ para $\vec{a}^t \vec{\beta}$

Dado o Modelo de Regressão Linear Múltipla e uma amostra, o intervalo a $(1-\alpha) \times 100\%$ de confiança para uma combinação linear dos parâmetros, $\vec{a}^t \vec{\beta} = a_0 \beta_0 + a_1 \beta_1 + ... + a_p \beta_p$, é:

$$\vec{\mathbf{a}}^t \vec{\mathbf{b}} - t_{\frac{\alpha}{2}[n-(p+1)]} \cdot \hat{\sigma}_{\vec{\mathbf{a}}^t \vec{\hat{\boldsymbol{\beta}}}} \quad , \quad \vec{\mathbf{a}}^t \vec{\mathbf{b}} + t_{\frac{\alpha}{2}[n-(p+1)]} \cdot \hat{\sigma}_{\vec{\mathbf{a}}^t \vec{\hat{\boldsymbol{\beta}}}} \quad \left[\right. ,$$

$$\mathbf{com} \quad \vec{\mathbf{a}}^t \vec{\mathbf{b}} = a_0 b_0 + a_1 b_1 + ... + a_p b_p \qquad \mathbf{e} \qquad \hat{\sigma}_{\vec{\mathbf{a}}^t \hat{\vec{\boldsymbol{\beta}}}} = \sqrt{QMRE \cdot \vec{\mathbf{a}}^t (\mathbf{X}^t \mathbf{X})^{-1} \vec{\mathbf{a}}}.$$

Fórmulas para a estimação de $\beta_i \pm \beta_i$

A fórmula geral $\hat{\sigma}_{\vec{a}^t \vec{\beta}} = \sqrt{QMRE \cdot \vec{a}^t (X^t X)^{-1} \vec{a}}$ admite uma expressão alternativa no caso particular duma soma ou diferença de dois β s.

Pela fórmula geral da variância duma soma ou diferença de v.a.s,

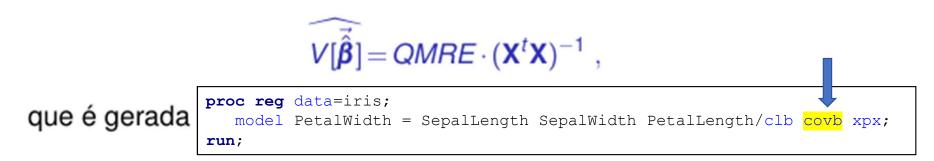
$$\begin{split} V[\hat{\beta}_{i} \pm \hat{\beta}_{j}] &= V[\hat{\beta}_{i}] + V[\hat{\beta}_{j}] \pm 2 \, Cov[\hat{\beta}_{i}, \hat{\beta}_{j}] \; . \\ \Leftrightarrow & \sigma_{\hat{\beta}_{i} \pm \hat{\beta}_{j}}^{2} &= \sigma^{2} \cdot \left[(\mathbf{X}^{t} \mathbf{X})_{[i+1,i+1]}^{-1} + (\mathbf{X}^{t} \mathbf{X})_{[j+1,j+1]}^{-1} \pm 2 \, (\mathbf{X}^{t} \mathbf{X})_{[i+1,j+1]}^{-1} \right] \; . \end{split}$$

Logo, o erro padrão de $\hat{\beta}_i \pm \hat{\beta}_j$ é:

$$\hat{\sigma}_{\hat{\beta}_i \pm \hat{\beta}_j} = \sqrt{QMRE \cdot \left[(\mathbf{X}^t \mathbf{X})_{[i+1,i+1]}^{-1} + (\mathbf{X}^t \mathbf{X})_{[j+1,j+1]}^{-1} \pm 2 (\mathbf{X}^t \mathbf{X})_{[i+1,j+1]}^{-1} \right]} \ .$$

ICs para combinações lineares

Numa RLM, o IC duma combinação linear genérica $\vec{a}^t \vec{\beta}$, precisa da matriz das (co)variâncias estimadas dos estimadores $\hat{\beta}$,



A matriz das (co)variâncias estimadas no exemplo RLM dos lírios é:

Covariance of Estimates								
Variable	Intercept	SepalLength	SepalWidth	PetalLength				
Intercept	0.0318157664	-0.005075942	-0.002486105	0.0015144174				
SepalLength	-0.005075942	0.0022568367	-0.001344002	-0.001065046				
SepalWidth	-0.002486105	-0.001344002	0.0023949317	0.000802941				
PetalLength	0.0015144174	-0.001065046	0.000802941	0.0005998259				

ICs para combinações lineares

O erro padrão estimado de $\hat{\beta}_1 + \hat{\beta}_2$

$$\hat{\sigma}_{\widehat{\beta}_{1}+\widehat{\beta}_{2}} = \sqrt{\widehat{V}[\widehat{\beta}_{1}] + \widehat{V}[\widehat{\beta}_{2}] + 2\widehat{Cov}[\widehat{\beta}_{1},\widehat{\beta}_{2}]}$$

$$\hat{\sigma}_{\widehat{\beta}_{1}+\widehat{\beta}_{2}} = \sqrt{0.0022568367 + 0.0023949317 + 2(-0.001344002)} = 0.04431439$$

Covariance of Estimates								
Variable	Intercept	SepalLength	SepalWidth	PetalLength				
Intercept	0.0318157664	-0.005075942	-0.002486105	0.0015144174				
SepalLength	-0.005075942	0.0022568367	-0.001344002	-0.001065046				
SepalWidth	-0.002486105	-0.001344002	0.0023949317	0.000802941				
PetalLength	0.0015144174	-0.001065046	0.000802941	0.0005998259				

Testes a combinações lineares dos parâmetros

Dado o Modelo de Regressão Linear Múltipla,

Testes de Hipóteses relativos a $\vec{a}^t \hat{\beta}$ Hipóteses: $H_0: \vec{\mathbf{a}}^t \vec{\boldsymbol{\beta}} \stackrel{\geq}{=} c \text{ vs. } H_1: \vec{\mathbf{a}}^t \vec{\boldsymbol{\beta}} \neq c$ Estatística do Teste: $T = \frac{\vec{a}^t \vec{\hat{\beta}} - \vec{a}^t \vec{\beta}|_{H_0}}{\hat{\sigma}_{zt\vec{\delta}}} \sim t_{n-(p+1)}$, se H_0 verdade Nível de significância do teste: α Região Crítica (Região de Rejeição): Rejeitar H₀ se $T_{calc} < -t_{\alpha[n-(p+1)]}$ (Unilateral esquerdo) (Bilateral) $|T_{calc}| > t_{\alpha/2[n-(p+1)]}$ $T_{calc} > t_{\alpha[n-(p+1)]}$ (Unilateral direito)

Intervalos de confiança para $E[Y|X_1=x_1,...,X_p=x_p]$

Como caso particular do resultado anterior, tem-se:

IC para o valor esperado de Y, dados os preditores

Dado o Modelo RLM e uma amostra com os valores $\vec{\mathbf{x}} = (x_1, x_2, ..., x_p)^t$ das variáveis preditoras, o valor esperado de Y,

$$\mu_{Y|\bar{x}} = E[Y|X_1 = x_1, ..., X_p = x_p] = \beta_0 + \beta_1 x_1 + ... + \beta_p x_p$$

é estimado por $\hat{\mu}_{Y|\vec{x}} = b_0 + b_1 x_1 + ... + b_p x_p$.

Um intervalo a $(1-\alpha) \times 100\%$ de confiança para $\mu_{Y|\bar{x}}$ é dado por:

$$\left] \quad \hat{\mu}_{\scriptscriptstyle Y|\vec{\mathbf{x}}} - t_{\scriptscriptstyle \underline{\alpha}} \left[_{\scriptscriptstyle n-(p+1)\right]} \cdot \hat{\sigma}_{\hat{\mu}_{\scriptscriptstyle Y|\vec{\mathbf{x}}}} \quad , \quad \hat{\mu}_{\scriptscriptstyle Y|\vec{\mathbf{x}}} + t_{\scriptscriptstyle \underline{\alpha}} \left[_{\scriptscriptstyle n-(p+1)\right]} \cdot \hat{\sigma}_{\hat{\mu}_{\scriptscriptstyle Y|\vec{\mathbf{x}}}} \quad \left[\right. ,$$

com
$$\hat{\sigma}_{\hat{\mu}_{Y|\vec{\mathbf{x}}}} = \sqrt{QMRE \cdot \vec{\mathbf{a}}^t(\mathbf{X}^t\mathbf{X})^{-1}\vec{\mathbf{a}}}$$
, onde $\vec{\mathbf{a}} = (1, x_1, x_2, ..., x_p)$.

Se p = 1, RLS

Fórmulas para uma regressão linear simples

Numa regressão linear simpes, a fórmula da variância de $\hat{\mu}_{Y|X}$ é:

$$\sigma_{\hat{\mu}_{Y|X}}^{2} = V[\hat{\mu}_{Y|X}] = \sigma^{2} \cdot \left[\frac{1}{n} + \frac{(x - \overline{x})^{2}}{\frac{(n-1) \cdot S_{X}^{2}}{2}} \right]$$

$$\implies \hat{\sigma}_{\hat{\mu}_{Y|X}}^{2} = QMRE \cdot \left[\frac{1}{n} + \frac{(x - \overline{x})^{2}}{\frac{(n-1) \cdot S_{X}^{2}}{2}} \right]$$

O intervalo de confiança para $\mu_{Y|X}$ na RLS é:

]
$$(b_0 + b_1 x) - t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\mu}_{Y|x}}$$
 , $(b_0 + b_1 x) + t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\hat{\mu}_{Y|x}}$ [.

A variabilidade duma observação individual de Y

Consideraram-se intervalos de confiança para o valor esperado de Y,

$$\mu_{Y|\vec{\mathbf{x}}} = E[Y|x_1 = x_1, x_2 = x_2, ..., x_p = x_p] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_p x_p ,$$

usam a variabilidade associada ao estimador $\hat{\mu}_{Y|\vec{x}}$:

$$\sigma_{\hat{\mu}_{Y|\vec{x}}}^2 = V[\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + ... + \hat{\beta}_p x_p] = \sigma^2 \cdot \vec{a}^t (X^t X)^{-1} \vec{a},$$

com $\vec{\mathbf{a}} = (1, x_1, x_2, ..., x_p).$

Uma observação individual de Y tem uma variabilidade adicional, pois:

$$Y = \mu_{Y|\vec{\mathbf{x}}} + \varepsilon = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_p x_p + \varepsilon.$$

A flutuação aleatória de observações individuais em torno do hiperplano é $V[\varepsilon] = \sigma^2$. Será necessário somar a variância associada à estimação do hiperplano e a variância das observações individuais:

$$\sigma_{\textit{Indiv}}^2 = V[\hat{\mu}_{Y|\vec{\mathbf{x}}}] + V[\varepsilon] = \sigma^2 \cdot \vec{\mathbf{a}}^t (\mathbf{X}^t \mathbf{X})^{-1} \vec{\mathbf{a}} + \sigma^2 = \sigma^2 \cdot \left[\vec{\mathbf{a}}^t (\mathbf{X}^t \mathbf{X})^{-1} \vec{\mathbf{a}} + 1 \right].$$

Intervalos de predição para Y

Podem obter-se intervalos de predição para uma observação individual de Y, associada aos valores $X_1 = x_1, ..., X_p = x_p$ das variáveis preditoras.

Nestes intervalos, a estimativa da variância duma observação individual de Y é a estimativa de σ_{Indiv}^2 , resultante de substituir σ^2 pelo QMRE amostral:

Intervalos de predição para observações individuais

$$\left] \quad \hat{\mu}_{Y|\vec{\mathbf{x}}} \ - \ t_{\frac{\alpha}{2}[n-(p+1)]} \cdot \hat{\sigma}_{indiv} \quad , \quad \hat{\mu}_{Y|\vec{\mathbf{x}}} \ + \ t_{\frac{\alpha}{2}[n-(p+1)]} \cdot \hat{\sigma}_{indiv} \quad \left[\right.$$

onde

$$\hat{\mu}_{Y|X} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p$$

е

$$\hat{\sigma}_{indiv} = \sqrt{QMRE \left[1 + \vec{\mathbf{a}}^{t}(\mathbf{X}^{t}\mathbf{X})^{-1}\vec{\mathbf{a}}\right]} \quad \text{com} \quad \vec{\mathbf{a}} = (1, x_1, x_2, ..., x_p).$$

Se p = 1, RLS

Fórmulas para a regressão linear simples

Na regressão linear simples usa-se a fórmula

$$\sigma_{Indiv}^{2} = \underbrace{\sigma^{2} \cdot \left[\frac{1}{n} + \frac{(x - \overline{x})^{2}}{\binom{n-1}{\cdot} S_{x}^{2}} \right]}_{=V[\hat{\mu}_{Y|\overline{x}}]} + \underbrace{\sigma^{2}}_{=V[\varepsilon]} = \sigma^{2} \cdot \left[1 + \frac{1}{n} + \frac{(x - \overline{x})^{2}}{\binom{n-1}{\cdot} S_{x}^{2}} \right].$$

Logo,

RLS: Intervalo de predição para observação individual de Y

Quer numa regressão linear simples, quer numa múltipla, estes intervalos são necessariamente de maior amplitude que os intervalos de confiança para $\mu_{Y|\bar{x}}$ (para igual nível de confiança $(1-\alpha) \times 100\%$).

Testando a qualidade do ajustamento global

Numa Regressão Linear, o modelo é inútil se fôr indistinguível do modelo nulo, i.e., do modelo de equação $Y_i = \beta_0 + \varepsilon_i$. O modelo nulo pode ser visto como um submodelo de qualquer modelo linear, em que todas as variáveis preditoras têm coeficiente nulo: $\beta_i = 0$, $\forall j > 0$.

O teste de ajustamento global visa testar se um dado modelo linear é significativamente diferente do modelo nulo.

As hipóteses em confronto são:

$$H_0: \beta_1 = \beta_2 = ... = \beta_p = 0$$

[MODELO COMPLETO \equiv MODELO NULO]
vs.
 $H_1: \exists j = 1,...,p$ t.q. $\beta_j \neq 0$
[MODELO COMPLETO $\not\equiv$ MODELO NULO]

NOTA: repare que β_0 não intervém nas hipóteses.

O teste de ajustamento global (cont.)

Definindo:

- O Quadrado Médio da Regressão como $QMR = \frac{SQR}{p}$.
- O Quadrado Médio Residual como $QMRE = \frac{SQRE}{n-(p+1)}$.

Sob a Hipótese Nula do teste de ajustamento global:

$$F = \frac{QMR}{QMRE} - F_{[p, n-(p+1)]}.$$

Esta é a estatística F do teste de ajustamento global.

Expressão alternativa para a estatística do teste F

A estatística do teste F de ajustamento global do modelo numa Regressão Linear Múltipla pode ser escrita na forma alternativa:

$$F = \frac{n-(p+1)}{p} \cdot \frac{R^2}{1-R^2} .$$

A estatística F é uma função crescente do coeficiente de determinação amostral R², o que justifica a natureza unilateral direita da região crítica.

As hipóteses do teste também se podem escrever como

$$H_0: \mathcal{R}^2 = 0$$
 vs. $H_1: \mathcal{R}^2 > 0$.

A hipótese $H_0: \mathcal{R}^2=0$ indica ausência de relação linear entre Y e o conjunto dos preditores. Corresponde a um ajustamento "péssimo" do modelo. A sua rejeição não garante um bom ajustamento.

O Teste *F* de ajustamento global do Modelo

Teste F de ajustamento global do modelo RLM

Hipóteses: H_0 : $\beta_1 = \beta_2 = ... = \beta_p = 0$ vs.

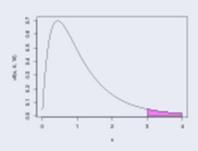
 H_1 : $\exists j = 1,...,p$ tal que $\beta_j \neq 0$.

Estatística do Teste: $F = \frac{QMR}{QMRE} \frown F_{[p,n-(p+1)]}$ se H_0 .

Nível de significância do teste: α

Região Crítica (Região de Rejeição): Unilateral direita

Rejeitar H_0 se $F_{calc} > f_{\alpha[p,n-(p+1)]}$



Outra formulação do teste F de ajustamento global

Teste F de ajustamento global do modelo RLM (alternativa)

```
Hipóteses: H_0: \mathcal{R}^2 = 0 vs. H_1: \mathcal{R}^2 > 0.
```

Estatística do Teste:
$$F = \frac{n-(p+1)}{p} \cdot \frac{R^2}{1-R^2} \frown F_{[p,n-(p+1)]}$$
 se H_0 .

Nível de significância do teste: α

Região Crítica (Região de Rejeição): Unilateral direita

Rejeitar
$$H_0$$
 se $F_{calc} > f_{\alpha(p,n-(p+1))}$

A hipótese nula $H_0: \mathcal{R}^2 = 0$ afirma que, na população, o coeficiente de determinação é nulo.

Informação Teste F de ajustamento Global

Analysis of Variance									
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F				
Model	3	81.18964	27.06321	734.39	<.0001				
Error	146	5.38030	0.03685						
Corrected Total	149	86.56993							

Root MSE	0.19197	R-Square	0.9379
Dependent Mean	1.19933	Adj R-Sq	0.9366
Coeff Var	16.00615		

O R^2 modificado (adjusted R^2)

O Coeficiente de Determinação usual define-se como:

$$R^2 = \frac{SQR}{SQT} = 1 - \frac{SQRE}{SQT}$$

O R^2 modificado, sendo $QMT = \frac{SQT}{n-1} = s_y^2$, é:

$$R_{mod}^2 = 1 - \frac{QMRE}{QMT} = 1 - \frac{SQRE}{SQT} \cdot \frac{n-1}{n-(p+1)} = 1 - (1-R^2) \cdot \frac{n-1}{n-(p+1)}.$$

Para qualquer modelo linear (com preditores), tem-se: $R_{mod}^2 < R^2$. Se $n \gg p+1$ (muito mais observações que parâmetros), $R^2 \approx R_{mod}^2$. Se n é pouco maior que p, $R_{mod}^2 \ll R^2$ (excepto se $R^2 \approx 1$).

 $\frac{QMRE}{QMT} = \frac{\hat{\sigma}^2}{s_y^2}$ é a proporção da variabilidade total de Y que permanece inexplicada após a introdução dos preditores. Logo, R_{mod}^2 é o ganho na explicação de s_y^2 associado ao modelo.

Root MSE	0.19197	R-Square	0.9379
Dependent Mean	1.19933	Adj R-Sq	0.9366
Coeff Var	16.00615		

O princípio da parcimónia na RLM

Recordemos o princípio da parcimónia na modelação: queremos um modelo que descreva adequadamente a relação entre as variáveis, mas que seja o mais simples (parcimonioso) possível.

Caso se disponha de um modelo de Regressão Linear Múltipla com um ajustamento considerado adequado, a aplicação deste princípio traduz-se em saber se será possível obter um modelo com menos variáveis preditoras, sem perder significativamente em termos de qualidade de ajustamento.

Modelo e Submodelos

Se dispomos de um modelo de Regressão Linear Múltipla, com relação de base

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5$$

chamamos submodelo a um modelo de regressão linear múltipla contendo apenas algumas das variáveis preditoras, e.g.,

$$Y = \beta_0 + \beta_2 x_2 + \beta_5 x_5 ,$$

Podemos identificar o submodelo pelo conjunto $\mathscr S$ das variáveis preditoras que pertencem ao submodelo. No exemplo, $\mathscr S=\{2,5\}$.

O modelo e o submodelo são idênticos se $\beta_j = 0$ para qualquer variável x_j cujo índice não pertença a \mathscr{S} .

Comparando modelo e submodelos

Para comparar um modelo e um seu submodelo (identificado pelo conjunto $\mathscr S$ dos índices das suas variáveis), precisamos de optar entre as hipóteses:

$$H_0: \beta_j = 0, \quad \forall j \notin \mathscr{S}$$
 vs. $H_1: \exists j \notin \mathscr{S}$ tal que $\beta_j \neq 0$. [SUBMODELO = MODELO]

NOTA: Esta discussão só envolve coeficientes β_j de variáveis preditoras (j > 0). O coeficiente β_0 faz sempre parte dos submodelos e não é relevante do ponto de vista da parcimónia.

Caso não se rejeite H_0 , opta-se pelo submodelo (mais parcimonioso).

Caso se rejeite H_0 , opta-se pelo modelo completo (ajusta-se significativamente melhor).

Este coeficiente β_0 não é relevante do ponto de vista da parcimónia: a sua presença não implica trabalho adicional de recolha de dados, nem de interpretação do modelo. Apenas permite um melhor ajustamento.

Estatística de teste para comparar modelo/submodelo

A estatística de teste compara as Somas de Quadrados Residuais do:

- modelo completo (referenciado pelo índice C); e do
- submodelo (referenciado pelo índice S)

Seja k o número de preditores do submodelo (k+1 parâmetros). Tem-se, sob H_0 ($\beta_i = 0$, para todas as variáveis x_i que não estão no submodelo):

$$F = \frac{\frac{SQRE_S - SQRE_C}{p - k}}{\frac{SQRE_C}{n - (p + 1)}} \qquad \qquad F_{[p - k, n - (p + 1)]},$$

Nota: Necessariamente $SQRE_s \ge SQRE_c$.

São os valores grandes da estatística que levantam dúvidas sobre H_0 .

O teste a um submodelo (teste F parcial)

Teste *F* de comparação dum modelo com um seu submodelo

Dado o Modelo de Regressão Linear Múltipla,

Hipóteses:

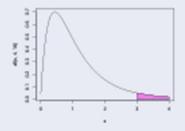
$$H_0: \beta_j = 0, \quad \forall j \notin \mathscr{S} \quad \text{vs.} \quad H_1: \exists j \notin \mathscr{S} \quad \text{tal que} \quad \beta_j \neq 0.$$

Estatística do Teste:
$$F = \frac{\frac{SQRE_S - SQRE_C}{p-k}}{\frac{SQRE_C}{n-(p+1)}} \qquad F_{[p-k,n-(p+1)]}, \text{ sob } H_0.$$

Nível de significância do teste: α

Região Crítica (Região de Rejeição): Unilateral direita

Rejeitar H_0 se $F_{calc} > f_{\alpha[p-k,n-(p+1)]}$



Expressão alternativa para a estatística do teste

A estatística do teste *F* parcial pode ser escrita na forma alternativa:

$$F = \frac{n - (p+1)}{p - k} \cdot \frac{R_C^2 - R_S^2}{1 - R_C^2}.$$

NOTA: A Soma de Quadrados Total apenas depende dos valores observados da variável resposta Y e não do modelo ajustado. Assim, SQT é igual no modelo completo e no submodelo.

As hipóteses do teste também se podem escrever como

$$H_0: \mathcal{R}_C^2 = \mathcal{R}_S^2$$
 vs. $H_1: \mathcal{R}_C^2 > \mathcal{R}_S^2$,

A hipótese H_0 indica que o grau de relacionamento linear entre Y e o conjunto dos preditores é idêntico no modelo e no submodelo.

Caso não se rejeite H_0 , opta-se pelo submodelo (mais parcimonioso). Caso se rejeite H_0 , opta-se pelo modelo completo (ajusta-se significativamente melhor).

Teste F parcial: formulação alternativa

Teste F de comparação dum modelo com um seu submodelo

Dado o Modelo de Regressão Linear Múltipla,

Hipóteses:

$$H_0: \mathscr{R}_C^2 = \mathscr{R}_S^2 \quad \text{vs.} \quad H_1: \mathscr{R}_C^2 > \mathscr{R}_S^2.$$

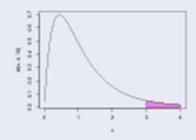
Estatística do Teste:

$$F = \frac{n-(p+1)}{p-k} \cdot \frac{R_C^2 - R_S^2}{1 - R_C^2} \qquad \qquad F_{[p-k, n-(p+1)]}, \text{ sob } H_0.$$

Nível de significância do teste: α

Região Crítica (Região de Rejeição): Unilateral direita

Rejeitar H_0 se $F_{calc} > f_{\alpha[p-k,n-(p+1)]}$



Relações dos testes *F* parcial

O teste de ajustamento global é equivalente a um teste *F* parcial comparando um modelo linear e o submodelo nulo (sem preditores).

Caso o modelo e submodelo difiram num único preditor, X_i , o teste F parcial é equivalente ao teste t com as hipóteses $H_0: \beta_j = 0$ vs. $H_1: \beta_j \neq 0$.

Nesse caso, não apenas as hipóteses dos dois testes são iguais, como a estatística do teste *F* parcial é o quadrado da estatística do teste *t* referido.

- as hipóteses dos dois testes são iguais $(H_0: \beta_i = 0 \text{ vs. } H_1: \beta_i \neq 0);$
- a estatística do teste F parcial é o quadrado da estatística do teste t referido:

$$F_{calc} = T_{calc}^2$$

Tem-se p - k = 1, e como é sabido, se uma variável aleatória T tem distribuição t_v , então o seu quadrado, T^2 tem distribuição $F_{1,v}$.

Numa regressão linear simples, o teste *t* ao declive da recta ser nulo é equivalente ao teste *F* de ajustamento global. A segunda destas estatística de teste é o quadrado da primeira.

Teste F Parcial de comparação de um modelo (com p preditores) com um seu submodelo (com k preditores)

Submodelo: PetalWidth = β_0 + β_3 PetalLength

Modelo completo: PetalWidth = β_0 + β_1 SepalLength+ β_2 SepalWidth+ β_3 PetalLength

$$H_0: \beta_1 = \beta_2 = 0$$
 vs. $H_1: \exists_{j=1,2}: \beta_j \neq 0$

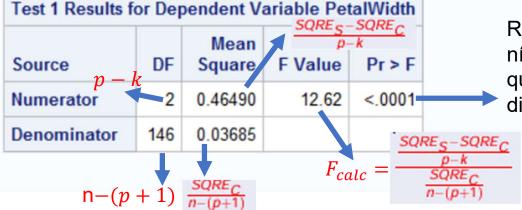
proc reg data=iris;

model PetalWidth = SepalLength SepalWidth PetalLength;
test SepalLength, SepalWidth;

run;

The SAS System

The REG Procedure Model: MODEL1



Rejeita-se a hipótese nula (para qualquer nível de significância usual), portanto, a qualidade do ajustamento dos dois modelos difere significativativamente.

Teste F Parcial de comparação de um modelo (com p preditores) com um seu submodelo (com k preditores)

Submodelo: PetalWidth = β_0 + β_3 PetalLength

Modelo completo: PetalWidth = β_0 + β_1 SepalLength+ β_2 SepalWidth+ β_3 PetalLength

De forma equivalente:

$$H_0: \mathscr{R}_C^2 = \mathscr{R}_S^2$$
 vs. $H_1: \mathscr{R}_C^2 > \mathscr{R}_S^2$

Os valores dos coeficientes de determinação amostrais ($R_s^2 = 0.9271 \, e \, R_c^2 = 0.9379$) são significativamente diferentes.

```
proc reg data=iris;
   model PetalWidth = PetalLength/clb covb xpx;
run;
```

Root MSE	0.20648	R-Square	0.9271
Dependent Mean	1.19933	Adj R-Sq	0.9266
Coeff Var	17.21659		

proc reg	data=in	ris;						
model	PetalWi	idth =	= Sep	pa]	LLeng	gth	Sepal	Width
PetalLeng	gth/clb	covb	хрх	R	CLI	CLM	;	
run;								

Root MSE	0.19197	R-Square	0.9379
Dependent Mean	1.19933	Adj R-Sq	0.9366
Coeff Var	16.00615		

Exercícios:

- 1) Usando os valores dos coeficientes de determinação dos dois modelos ajustados verifique que o valor do $F_{ParcialCalc} = 12.62$ (slide anterior);
- 2) Usando os valores das somas dos quadrados dos resíduos dos dois modelos ajustados verifique que o valor do $F_{ParcialCalc} = 12.62$ (slide anterior).

Exercícios (continuação):

Relembrando alguns resultados do ajustamento do modelo completo:

```
proc reg data=iris;
   model PetalWidth = SepalLength SepalWidth PetalLength/clb covb xpx R CLI CLM ;
run;
```

			Parameter	Estimate	S		
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confide	nce Limits
Intercept	1	-0.24031	0.17837	-1.35	0.1800	-0.59283	0.11221
SepalLength	1	-0.20727	0.04751	-4.36	<.0001	-0.30115	-0.11338
SepalWidth	1	0.22283	0.04894	4.55	<.0001	0.12611	0.31955
PetalLength	1	0.52408	0.02449	21.40	<.0001	0.47568	0.57249

- 3) a) Utilize um teste F parcial para ver se é possível concluir que os modelos com e sem o preditor SepalLength têm ajustamento significativamente diferente (utilize $\alpha = 0.05$).
 - b) Qual o coeficiente de determinação do submodelo resultante da exclusão dessa variável?

Como escolher um submodelo?

O teste *F* parcial (teste aos modelos encaixados) permite-nos optar entre um modelo e um seu submodelo. Por vezes, um submodelo pode ser sugerido por:

- razões de índole teórica, sugerindo que determinadas variáveis preditoras não sejam, na realidade, importantes para influenciar os valores de Y.
- razões de índole prática, como a dificuldade, custo ou volume de trabalho associado à recolha de observações para determinadas variáveis preditoras.

Nestes casos, pode ser claro que submodelo(s) se deseja testar.

Como escolher um submodelo? (cont.)

Mas em muitas situações não é evidente qual o subconjunto de variáveis preditoras que se deseja considerar no submodelo. Pretende-se apenas ver se o modelo é simplificável. Nestes casos, a opção por um submodelo não é um problema fácil.

Dadas p variáveis preditoras, o número de subconjuntos, de qualquer cardinalidade, excepto 0 (modelo nulo) e p (o modelo completo) que é possível escolher é dado por $2^p - 2$. A tabela seguinte indica o número desses subconjuntos para p = 5, 10, 15, 20, 30.

р	$2^{p}-2$
5	30
10	1 022
15	32 766
20	1 048 574
30	1 073 741 822

Para valores de p pequenos, é possível analisar todos os possíveis subconjuntos. Com algoritmos e rotinas informáticas adequadas, a pesquisa completa de todos os possíveis subconjuntos ainda é possível para valores grandes de p (até $p \approx 35$). Mas para p muito grande, uma pesquisa completa é computacionalmente inviável.

Não é legítimo optar pela exclusão de várias variáveis preditoras em simultâneo, com base nos testes t à significância de cada coeficiente β_i no modelo completo.

De facto, os testes t aos coeficientes β_j admitem que todas as restantes variáveis pertencem ao modelo. A exclusão de um qualquer preditor altera o ajustamento: altera os valores estimados b_j e os respectivos erros padrão das variáveis que permanecem no submodelo. Pode acontecer que um preditor seja dispensável num modelo completo, mas deixe de o ser num submodelo, ou viceversa.

Algoritmos de pesquisa sequenciais

Caso não esteja disponível *software* apropriado, ou se o número *p* de preditores for demasiado grande, pode recorrer-se a algoritmos de pesquisa que simplificam uma regressão linear múltipla sem analisar todo os possíveis submodelos e sem a garantia de obter os melhores subconjuntos.

Vamos considerar um algoritmo que, em cada passo, exclui uma variável preditora, até alcançar uma condição de paragem considerada adequada, ou seja, um algoritmo de exclusão sequencial (backward elimination).

Existem variantes deste algoritmo, não estudadas aqui:

- algoritmo de inclusão sequencial (forward selection).
- algoritmos de exclusão/inclusão alternada (stepwise selection).

O algoritmo de exclusão sequencial com testes aos β_j

- ajustar o modelo completo, com os p preditores;
- ② definir um nível de significância α para os testes de hipóteses a $\beta_i = 0$;
- opara todas as variáveis rejeita-se $H_0: \beta_i = 0$?
 - Se sim: não é possível simplificar o modelo (passar ao ponto 4).
 - Se não: variáveis em que não se rejeita H₀ são dispensáveis (candidatas à exclusão).
 - se apenas existe uma candidata a sair, excluir essa variável;
 - se existir mais do que uma variável candidata a sair, excluir a variável associada ao maior p-value (isto é, ao valor da estatística t mais próxima de zero)

Reajustar o modelo após a exclusão da variável e repetir este ponto 3

Quando não existirem variáveis candidatas a sair, ou quando sobrar um único preditor, o algoritmo pára. Tem-se então o submodelo final.

Critério de Informação de Akaike

O R disponibiliza funções para automatizar pesquisas sequenciais de submodelos, semelhantes à que aqui foi enunciada, mas em que o critério de exclusão duma variável em cada passo se baseia no Critério de Informação de Akaike (AIC).

Critério de Informação de Akaike (AIC)

O AIC é uma medida geral da qualidade de ajustamento de modelos. No contexto duma Regressão Linear Múltipla com *k* variáveis preditoras, define-se como

$$AIC = n \cdot \ln \left(\frac{SQRE_k}{n} \right) + 2(k+1).$$

Nota: O AIC pode tomar valores negativos.

Interpretando o AIC

$$AIC = n \cdot \ln \left(\frac{SQRE_k}{n} \right) + 2(k+1)$$

- a primeira parcela é função crescente de SQRE_k, i.e., quanto melhor o ajustamento, mais pequena a primeira parcela;
- a segunda parcela mede a complexidade do modelo (k+1 é o número de parâmetros), pelo que quanto mais parcimonioso o modelo, mais pequena a segunda parcela.

Assim, o AIC depende simultaneamente da qualidade do ajustamento e da simplicidade do modelo.

Um modelo para a variável resposta Y é considerado melhor que outro se tiver um AIC menor (quando ajustados com os mesmos dados).

Algoritmo de exclusão sequencial com base no AIC

Pode definir-se um algoritmo de exclusão sequencial, com base no critério AIC:

- ajustar o modelo completo e calcular o respectivo AIC.
- ajustar cada submodelo com menos uma variável e calcular o respectivo AIC.
- Se nenhum dos AICs dos submodelos considerados for inferior ao AIC do modelo anterior, o algoritmo termina sendo o modelo anterior o modelo final.
 - Caso alguma das exclusões reduza o AIC, efectua-se a exclusão que mais reduz o AIC e regressa-se ao ponto anterior.

As duas variantes dos algoritmos

Os algoritmos de exclusão sequencial baseados nos testes *t* ou no AIC coincidem nas variáveis a excluir, podendo diferir apenas no momento de paragem.

Em geral, um algoritmo de exclusão sequencial baseado no AIC é mais cauteloso na exclusão, sobretudo se o valor de α usado nos testes t for baixo. Nos algoritmos baseados nos testes t, é aconselhável usar valores mais elevados de α , como $\alpha = 0.10$.

Um algoritmo de exclusão sequencial não garante a identificação do "melhor submodelo" com um dado número de preditores. Apenas identifica, de forma computacionalmente ligeira, submodelos "bons".

Deve ser usado com bom senso e o submodelo obtido cruzado com outras considerações (e.g., o custo ou dificuldade de obtenção de cada variável, ou o papel que a teoria relativa ao problema em questão reserva a cada preditor).

Exemplo: prever a percentagem de músculo em carcaça de porcos a partir de 7 preditores

```
proc reg data=porcos;
model Musculo = Area Gordurasubcut Peso Rendimento Gordurarenalpel Comprimento LarguraAnca /clb
covb xpx R CLI CLM;
output out=out_reg p=predicted_value;
test Area, Gordurasubcut, Peso, Rendimento;
RUN;
quit;
```

		P	arameter E	stimates			
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confidence Limi	
Intercept	1	54.11487	9.27061	5.84	<.0001	33.91594	74.31380
Area	1	0.06200	0.70162	0.09	0.9310	-1.46670	1.59070
Gordurasubcut	1	-0.93861	0.36030	-2.61	0.0230	-1.72363	-0.15359
Peso	1	0.24489	0.26196	0.93	0.3683	-0.32587	0.81565
Rendimento	1	0.00623	0.08323	0.07	0.9416	-0.17511	0.18756
Gordurarenalpel	1	-0.01436	0.00714	-2.01	0.0673	-0.02991	0.00119
Comprimento	1	0.01774	0.04832	0.37	0.7199	-0.08755	0.12302
LarguraAnca	1	0.11974	0.06255	1.91	0.0797	-0.01654	0.25602

- Há 6 preditores cuja exclusão individual é admissível (com $\alpha = 0.05$).
- Mas não é legítimo concluir que Area, Peso, Rendimento, Gordurarenalpel, Comprimento e LarguraAnca são dispensváveis em conjunto.

O algoritmo de exclusão sequencial com testes aos β_j

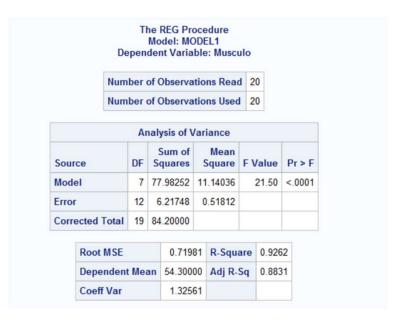
 $(\mathsf{com}\ \alpha=0.10)$

ajustar o modelo completo, com os p preditores

Modelo inicial, Completo, p = 7

model Musculo = Area Gordurasubcut Peso Rendimento
Gordurarenalpel Comprimento LarguraAnca

Parameter Estimates										
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confide	nce Limits			
Intercept	1	54.11487	9.27061	5.84	<.0001	33.91594	74.31380			
Агеа	1	0.06200	0.70162	0.09	0.9310	-1.46670	1.59070			
Gordurasubcut	1	-0.93861	0.36030	-2.61	0.0230	-1.72363	-0.15359			
Peso	1	0.24489	0.26196	0.93	0.3683	-0.32587	0.81565			
Rendimento	1	0.00623	0.08323	0.07	0.9416	-0.17511	0.18756			
Gordurarenalpel	1	-0.01436	0.00714	-2.01	0.0673	-0.02991	0.00119			
Comprimento	1	0.01774	0.04832	0.37	0.7199	-0.08755	0.12302			
LarguraAnca	1	0.11974	0.06255	1.91	0.0797	-0.01654	0.25602			



para todas as variáveis rejeita-se $H_0: \beta_j = 0$? $(\alpha = 0.10)$

se existir mais do que uma variável candidata a sair, excluir a variável associada ao maior *p-value* (isto é, ao valor da estatística *t* mais próxima de zero)

Reajustar o modelo após a exclusão da variável rendimento

model Musculo = Area Gordurasubcut Peso Gordurarenalpel Comprimento LarguraAnca

	Parameter Estimates										
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confidence Limits					
Intercept	1	54.50052	7.40466	7.36	<.0001	38.50373	70.49732				
Area	1	0.05151	0.66065	0.08	0.9390	1.37575	1.47877				
Gordurasubcut	1	-0.94880	0.32059	-2.96	0.0111	-1.64138	-0.25621				
Peso	1	0.25275	0.23057	1.10	0.2929	-0.24536	0.75087				
Gordurarenalpel	1	-0.01427	0.00677	-2.11	0.0549	-0.02889	0.00034730				
Comprimento	1	0.01672	0.04456	0.38	0.7135	-0.07955	0.11300				
LarguraAnca	1	0.11927	0.05981	1.99	0.0675	-0.00994	0.24849				

para todas as variáveis rejeita-se $H_0: \beta_j = 0$?

 $(\alpha=0.10)$

se existir mais do que uma variável candidata a sair, excluir a variável associada ao maior *p-value* (isto é, ao valor da estatística *t* mais próxima de zero)

Reajustar o modelo após a exclusão da variável Area

model Musculo = Gordurasubcut Peso Gordurarenalpel Comprimento LarguraAnca

Parameter Estimates									
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confide	nce Limits		
Intercept	1	54.95081	4.46691	12.30	<.0001	45.37024	64.53139		
Gordurasubcut	1	-0.96660	0.21686	-4.46	0.0005	-1.43173	-0.50147		
Peso	1	0.25637	0.21768	1.18	0.2585	-0.21050	0.72325		
Gordurarenalpel	1	-0.01457	0.00534	-2.73	0.0162	-0.02602	-0.00313		
Comprimento	1	0.01747	0.04195	0.42	0.6834	-0.07251	0.10745		
LarguraAnca	1	0.11937	0.05764	2.07	0.0573	-0.00425	0.24298		

para todas as variáveis rejeita-se $H_0: \beta_j = 0$? $(\alpha = 0.10)$

se existir mais do que uma variável candidata a sair, excluir a variável associada ao maior *p-value* (isto é, ao valor da estatística *t* mais próxima de zero)

Reajustar o modelo após a exclusão da variável Comprimento

model Musculo = Gordurasubcut Peso Gordurarenalpel LarguraAnca

	Parameter Estimates									
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confide	ence Limits			
Intercept	1	55.95724	3.65149	15.32	<.0001	48.17427	63.74020			
Gordurasubcut	1	-0.98841	0.20456	-4.83	0.0002	-1.42443	-0.55240			
Peso	1	0.26808	0.20982	1.28	0.2208	-0.17915	0.71531			
Gordurarenalpel	1	-0.01423	0.00512	-2.78	0.0141	-0.02515	-0.00331			
LarguraAnca	1	0.12268	0.05549	2.21	0.0430	0.00441	0.24095			

para todas as variáveis rejeita-se $H_0: \beta_j = 0$?

 $(\alpha = 0.10)$

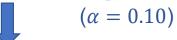
se apenas existe uma candidata a sair, excluir essa variável;

Reajustar o modelo após a exclusão da variável Peso

model Musculo = Gordurasubcut Gordurarenalpel LarguraAnca

Parameter Estimates										
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confide	% Confidence Limits			
Intercept	1	60.21488	1.52200	39.56	<.0001	56.98839	63.44138			
Gordurasubcut	1	-0.92545	0.20242	-4.57	0.0003	-1.35456	-0.49633			
Gordurarenalpel	1	-0.01721	0.00465	-3.70	0.0019	-0.02707	-0.00734			
LarguraAnca	1	0.11380	0.05613	2.03	0.0596	-0.00519	0.23279			

para todas as variáveis rejeita-se $H_0: \beta_j = 0$?



Quando não existirem variáveis candidatas a sair, ou quando sobrar um único preditor, o algoritmo pára. Tem-se então o submodelo final.

Submodelo final:

model Musculo = Gordurasubcut Gordurarenalpel LarguraAnca

Root MSE	0.66078	R-Square	0.9170
Dependent Mean	54.30000	Adj R-Sq	0.9015
Coeff Var	1.21690		

Parameter Estimates										
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confide	ence Limits			
Intercept	1	60.21488	1.52200	39.56	<.0001	56.98839	63.44138			
Gordurasubcut	1	-0.92545	0.20242	-4.57	0.0003	-1.35456	-0.49633			
Gordurarenalpel	1	-0.01721	0.00465	-3.70	0.0019	-0.02707	-0.00734			
LarguraAnca	1	0.11380	0.05613	2.03	0.0596	-0.00519	0.23279			

Algoritmo de exclusão sequencial com base no AIC

Um modelo para a variável resposta Y é considerado melhor que outro se tiver um AIC menor (quando ajustados com os mesmos dados).

proc reg data=porcos;

model Musculo = Area Gordurasubcut Peso

Rendimento Gordurarenalpel Comprimento

LarguraAnca /clb covb xpx R CLI CLM

selection=adjrsq aic;

RUN;

Number in Model	Adjusted R-Square	R-Square	AIC	Variables in Model
4	0.9052	0.9252	-13.1025	Gordurasubcut Peso Gordurarenalpel LarguraAnca
3	0.9015	0.9170	-13.0365	Gordurasubcut Gordurarenalpel LarguraAnca
5	0.8997	0.9261	-11.3487	Gordurasubcut Peso Gordurarenalpel Comprimento LarguraAnca
5	0.8987	0.9253	-11.1425	Area Gordurasubcut Peso Gordurarenalpel LarguraAnca
5	0.8985	0.9252	-11.1101	Gordurasubcut Peso Rendimento Gordurarenalpel LarguraAnca
4	0.8971	0.9188	-11.4592	Gordurasubcut Gordurarenalpel Comprimento LarguraAnca
4	0.8962	0.9181	-11.2878	Area Gordurasubcut Gordurarenalpel LarguraAnca
4	0.8956	0.9176	-11.1738	Gordurasubcut Rendimento Gordurarenalpel LarguraAnca
6	0.8920	0.9261	-9.3580	Area Gordurasubcut Peso Gordurarenalpel Comprimento LarguraAnca
6	0.8920	0.9261	-9.3543	Gordurasubcut Peso Rendimento Gordurarenalpel Comprimento LarguraAnca
5	0.8915	0.9201	-9.7781	Gordurasubcut Rendimento Gordurarenalpel Comprimento LarguraAnca
6	0.8909	0.9253	-9.1440	Area Gordurasubcut Peso Rendimento Gordurarenalpel LarguraAnca
5	0.8905	0.9193	-9.5898	Area Gordurasubcut Gordurarenalpel Comprimento LarguraAnca
5	0.8900	0.9189	-9.5004	Area Gordurasubcut Rendimento Gordurarenalpel LarguraAnca
6	0.8842	0.9208	-7.9614	Area Gordurasubcut Rendimento Gordurarenalpel Comprimento LarguraAnca
2	0.8834	0.8957	-10.4631	Gordurasubcut Gordurarenalpel

150

A Validação do Modelo (análise dos resíduos)

TODA a inferência feita até aqui admitiu a validade do Modelo Linear, e em particular, dos pressupostos relativos aos erros aleatórios: Normais, de média zero, variância homogénea e independentes.

Uma análise de regressão não fica completa sem que haja uma validação dos pressupostos do modelo.

A validação dos pressupostos relativos aos erros aleatórios (que são desconhecidos) faz-se através dos seus preditores, os resíduos.

A análise de Resíduos e outros diagnósticos

Uma análise de regressão linear não fica completa sem o estudo dos resíduos e de alguns outros diagnósticos.

O modelo linear admite que $\varepsilon_i \sim \mathcal{N}(0, \sigma^2) \quad \forall i = 1, ..., n$.

Sob o modelo linear, os resíduos têm a seguinte distribuição:

$$E_i \sim \mathcal{N}\left(0, \sigma^2(1-h_{ii})\right) \quad \forall i=1,...,n,$$

sendo h_{ii} o i-ésimo elemento diagonal da matriz $\mathbf{H} = \mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t$ de projecção ortogonal sobre o subespaço $\mathscr{C}(\mathbf{X})$.

Este resultado demonstra-se mais facilmente considerando o vector dos resíduos, $\vec{E} = \vec{Y} - \hat{\vec{Y}} = \vec{Y} - H\vec{Y} = (I_n - H)\vec{Y}$.

Propriedades dos Resíduos sob o modelo linear

Teorema (Distribuição dos Resíduos no Modelo Linear)

Dado o Modelo Linear, tem-se:

$$\vec{\mathbf{E}} \sim \mathcal{N}_n \left(\vec{\mathbf{0}}, \sigma^2 (\mathbf{I}_n - \mathbf{H}) \right)$$
 sendo $\vec{\mathbf{E}} = (\mathbf{I}_n - \mathbf{H}) \vec{\mathbf{Y}}$.

Como no Modelo Linear $\vec{Y} \sim \mathcal{N}(X\vec{\beta}, \sigma^2I_n)$, o vector dos resíduos $\vec{E} = (I_n - H)\vec{Y}$, tem distribuição Multinormal em sentido generalizado

O vector esperado de **E** resulta das propriedades do acetato 125:

- $E[\vec{E}] = E[(I_n H)\vec{Y}] = (I_n H)E[\vec{Y}] = (I_n H)X\vec{\beta} = \vec{0}$, pois $X\vec{\beta} \in \mathcal{C}(X)$, logo permanece invariante sob a projecção: $HX\vec{\beta} = X\vec{\beta}$.
- Pelas propriedades do acetato 126 e o facto de H ser simétrica (H^t = H) e idempotente (HH = H), tem-se:
 V[Ĕ] = V[(I_n H)Ÿ] = (I_n H)V[Ÿ](I_n H)^t = σ²·(I_n H).

Propriedades dos Resíduos no Modelo Linear (cont.)

Nota: Embora no modelo RL os erros aleatórios sejam independentes, os resíduos não são variáveis aleatórias independentes, pois as covariâncias entre resíduos diferentes são (em geral), não nulas:

$$cov(E_i, E_j) = -\sigma^2 \cdot h_{ij}$$
, se $i \neq j$,

onde h_{ij} indica o elemento da linha i e coluna j da matriz \mathbf{H} .

Se $\vec{\mathbf{E}} \frown \mathscr{N}_n \Big(\vec{\mathbf{0}}, \sigma^2 (\mathbf{I}_n - \mathbf{H}) \Big)$, então cada resíduo tem distribuição:

$$E_i \sim \mathcal{N}\left(0, \sigma^2(1-h_{ii})\right),$$

onde h_{ii} é o i-ésimo elemento diagonal de **H** e

$$\frac{E_i}{\sqrt{\sigma^2(1-h_{ii})}} \sim \mathcal{N}(0,1).$$

Resíduos habituais : $E_i = Y_i - \hat{Y}_i$;

Resíduos (internamente) estandardizados : $R_i = \frac{E_i}{\sqrt{QMRE \cdot (1-h_{ii})}}$.

Resíduos Studentizados (ou externamente estandardizados):

$$T_i = \frac{E_i}{\sqrt{QMRE_{[-i]} \cdot (1 - h_{ii})}}$$

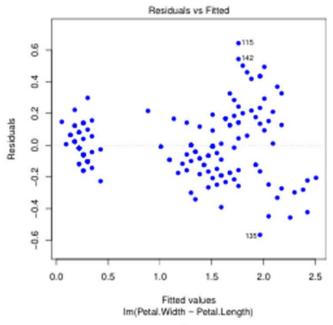
sendo $QMRE_{[-i]}$ o valor de QMRE resultante de um ajustamento da Regressão excluíndo a i-ésima observação (associada ao resíduo E_i).

Como
$$\frac{E_i}{\sqrt{\sigma^2(1-h_{ii})}} \sim \mathcal{N}(0,1)$$
, definem-se resíduos normalizados:

Para grandes amostras, os R_i são aproximadamente $\mathcal{N}(0,1)$.

Nas regressões lineares, avalia-se a validade dos pressupostos do modelo através de gráficos de resíduos. Não se efectuam testes de Normalidade, já que os resíduos não são (em geral) independentes.

Validação do modelo: (1) Gráficos de resíduos vs. \hat{Y}_i Gráfico indispensável: Resíduos (usuais) vs. Valores ajustados de Y.



- Os resíduos devem estar aproximadamente numa banda horizontal em torno de zero.
- Não deve existir qualquer padrão aparente. Sendo válido o Modelo RL, cor(E_i, Ŷ_i) = 0.

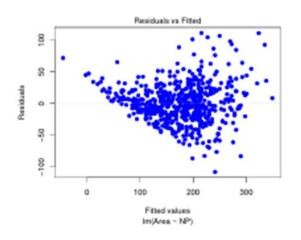
Possíveis padrões indicativos de problemas

Num gráfico de E_i vs. \hat{Y}_i podem surgir padrões problemáticos:

Curvatura na disposição dos resíduos Indica violação da hipótese de linearidade entre y e os preditores.

Gráfico em forma de funil Indica violação da hipótese de homogeneidade de variâncias.

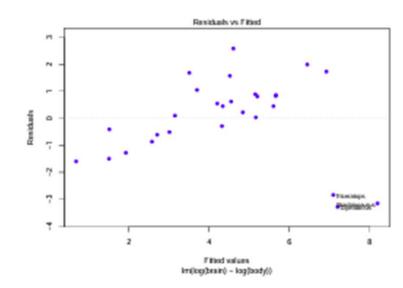
Um ou mais resíduos muito destacados Indica a existência de observações atípicas.



Um exemplo de resíduos em forma de funil, e sugerindo alguma curvatura na relação entre as duas variáveis

Padrões indicativos de problemas (cont.)

Um ou mais resíduos muito destacados e/ou banda oblíqua: Indica possíveis observações atípicas.



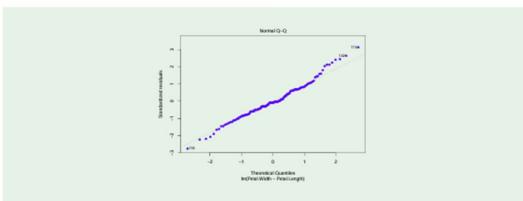
Validação do modelo: (2) Gráficos para avaliar a Normalidade

para grandes amostras os resíduos estandardizados $R_i = \frac{E_i}{\sqrt{QMRE \cdot (1-h_{ii})}}$, devem ser aproximadamente $\mathcal{N}(0,1)$

O pressuposto de erros aleatórios Normais pode ser validado com:

• um qq-plot que confronte os quantis empíricos dos n resíduos standardizados, com os quantis teóricos numa $\mathcal{N}(0,1)$.

Um qq-plot concordante com a hipótese de Normalidade dos erros aleatórios deverá apresentar colinearidade aproximada. O exemplo seguinte sugere algum desvio à Normalidade para os resíduos mais extremos.



Este *qq-plot* sugere algum desvio para os resíduos mais extremos, mas não em quantidade ou de forma suficientemente severa para pôr em dúvida o pressuposto da Normalidade dos erros aleatórios.

Validação do modelo: (3) Gráficos para avaliar independência

Dependência entre erros aleatórios pode surgir como resultado de:

- correlação cronológica;
- correlação espacial.

Nesse caso, pode ser útil inspeccionar gráficos de resíduos vs. ordem de observação ou distribuição espacial dos resíduos, para verificar se existem padrões que sugiram falta de independência. Nesse caso, modelos alternativos para series temporais ou dados espaciais podem ser necessários.

Validação do modelo: (4) Gráficos de resíduos vs. preditores

A presença de não-linearidade em gráficos de resíduos vs. preditores individuais pode sugerir a necessidade de transformações desses preditores.

Observações atípicas

Outras ferramentas de diagnóstico visam identificar observações individuais que merecem ulterior análise.

Observações atípicas (*outliers* in English). Conceito sem definição rigorosa, procura designar observações que se distanciam da relação linear de fundo entre *Y* e as variáveis preditoras.

Muitas vezes surgem associadas a resíduos grandes (em módulo). Em particular, e como os resíduos Studentizados têm distribuição aproximadamente $\mathcal{N}(0,1)$ para n grande, observações para as quais $|R_i| > 3$ (ou $|T_i| > 3$) podem ser classificadas como atípicas.

Mas por vezes, observações distantes da tendência geral podem afectar o próprio ajustamento do modelo, e não serem facilmente identificáveis a partir dos seus resíduos.

As chamadas "observações alavanca"

Define-se o valor do efeito alavanca (leverage) da i-ésima observação como sendo o i-ésimo valor diagonal da matriz \mathbf{H} : $h_{ii} = \mathbf{H}_{(i,i)}$.

Como $\hat{\mathbf{Y}} = \mathbf{H}\hat{\mathbf{Y}}$, tem-se $\hat{y}_i = \sum_{j=1}^n h_{ij} y_j$ (cada valor ajustado é combinação linear dos valores observados). O efeito alavanca h_{ii} é a ponderação associada a y_i na definição do valor ajustado correspondente, \hat{y}_i . Não deveria ser excessivo.

Observações alavanca (*leverage points*) são observações com *h_{ii}* elevado, que tendem a "atrair" a hipersuperfície ajustada numa regressão.

Como $V[E_i] = \sigma^2 (1 - h_{ii})$, se h_{ii} é elevado, a variância do resíduo E_i é baixa e o resíduo tende a estar perto da sua média (zero). Ou seja, a superfície ajustada tende a passar próximo desse ponto.

Observações alavanca (cont.)

Verifica-se, para qualquer observação:

$$\frac{1}{n} \leq h_{ii} \leq 1.$$

O valor médio das observações alavanca numa regressão linear é a razão entre o no. de parâmetros e o no. de observações:

$$\overline{h} = \frac{p+1}{n}$$
,

Logo, quanto mais observações, menor o efeito alavanca médio.

Observações alavanca (cont.)

Observações com um efeito alavanca elevado podem, ou não, estar dispostas com a mesma tendência de fundo que as restantes observações, i.e., podem, ou não, ser atípicas (outliers).

Efeito alavanca numa regressão linear Simples

Numa regressão linear simples, tem-se

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \overline{x})^2}{\frac{(n-1)\cdot S_x^2}{2}}.$$

Assim, numa RLS, o efeito alavanca da observação i depende do valor x_i em relação à média \overline{x} : quanto maior $(x_i - \overline{x})^2$, maior h_{ii} . O maior efeito alavanca tem de pertencer a uma das duas observações mais extrema em x.

Numa regressão linear múltipla, os maiores efeitos alavanca correspondem às observações em que os valores dos preditores estão mais afastados do vector das médias dos preditores.

Observações influentes

Observações influentes são observações que, se retiradas da análise, gerariam variações assinaláveis no conjunto dos valores ajustados de Y e nos parâmetros ajustados, b_i .

Medida de influência frequente é a distância de Cook, definida como:

$$D_{i} = \frac{\sum_{j=1}^{n} (\hat{y}_{j} - \hat{y}_{[-i]_{j}})^{2}}{(p+1) \cdot QMRE},$$

sendo $\hat{y}_{[-i]_i}$ o valor ajustado da observação i, obtido estimando os β_j s sem a observação i. Expressão equivalente é:

$$D_i = R_i^2 \cdot \left(\frac{h_{ii}}{1 - h_{ii}}\right) \cdot \frac{1}{p + 1}$$

Quanto maior D_i , maior é a influência da i-ésima observação.

É frequente considerar $D_i > 0.5$ como limiar de observação influente.

Uma prevenção

Observações atípicas, influentes ou alavanca não são o mesmo conceito, embora possam estar relacionados.

$$D_i = R_i^2 \cdot \left(\frac{h_{ii}}{1 - h_{ii}}\right) \cdot \frac{1}{p + 1}$$

 R_i^2 grande e h_{ii} grande $\Rightarrow D_i$ grande (observação influente)

 R_i^2 pequeno e h_{ii} pequeno $\Rightarrow D_i$ pequeno (observação não influente)

 R_i^2 grande e h_{ii} pequeno (ou viceversa) – D_i pode, ou não, ser grande

(Se obs. i é, ou não, influente depende da grandeza relativa de R_i^2 e h_{ii})

Estes diagnósticos servem sobretudo para identificar observações que merecem maior atenção e consideração.

ESTUDOS DOS RESÍDUOS

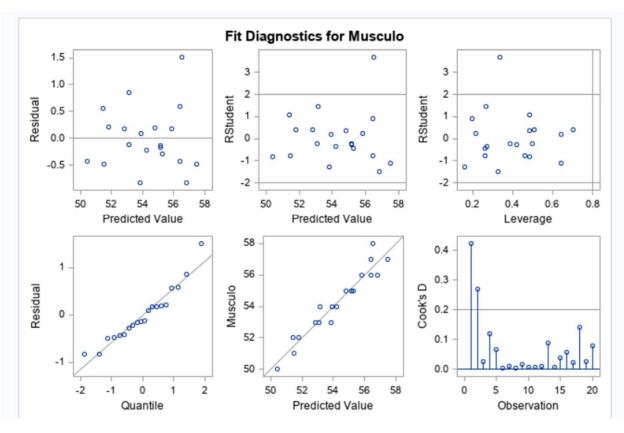
proc reg data=porcos plots=diagnostics; model Musculo = Area Gordurasubcut Peso Rendimento Gordurarenalpel Comprimento LarguraAnca /clb covb xpx R CLI CLM R P ; output out=diagnostics r=r student=int_r rstudent=ext_r h=leverage cookd=cooksd p=predicted; Resíduos usuais, Ei Resíduos Resíduos RUN; Distâncias **Efeitos** standardizados, Ri studentizados, Ti de Cook, Di alavanca, hii proc print data=diagnostics;

Sum of Residuals	0
Sum of Squared Residuals	6.30036
Predicted Residual SS (PRESS)	10.26395

Exercícios

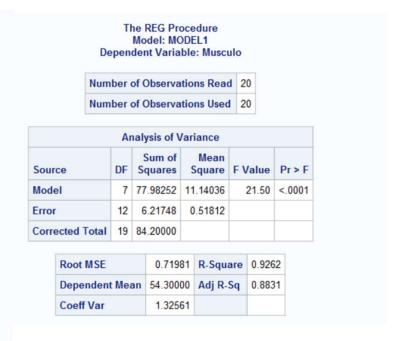
run;

a) Estude os gráficos de resíduos e outros diagnósticos



b) Os resultados do ajustamento do modelo com todos os preditores apresentam-se seguidamente:

Parameter Estimates										
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confidence Limit				
Intercept	1	54.11487	9.27061	5.84	<.0001	33.91594	74.31380			
Area	1	0.06200	0.70162	0.09	0.9310	-1.46670	1.59070			
Gordurasubcut	1	-0.93861	0.36030	-2.61	0.0230	-1.72363	-0.15359			
Peso	1	0.24489	0.26196	0.93	0.3683	-0.32587	0.81565			
Rendimento	1	0.00623	0.08323	0.07	0.9416	-0.17511	0.18756			
Gordurarenalpel	1	-0.01436	0.00714	-2.01	0.0673	-0.02991	0.00119			
Comprimento	1	0.01774	0.04832	0.37	0.7199	-0.08755	0.12302			
LarguraAnca	1	0.11974	0.06255	1.91	0.0797	-0.01654	0.25602			



					Output	t Statistic	s						
Obs	Dependent Variable	Predicted Value	Std Error Mean Predict			95% CL Predict		Resíduos usuais Residual	usuais standa Std Error Stu		zados Cook, Di		
1	58	56.4895	0.4179	55.5789	57.4001	54.6760	58.3030	1.5105	0.586	2.577	0.422		
2	57	57.4744	0.5758	56.2198	58.7290	55.4660	59.4828	-0.4744	0.432	-1.098	0.268		
3	57	56.4076	0.3174	55.7161	57.0992	54.6936	58.1217	0.5924	0.646	0.917	0.025		
4	56	56.8284	0.4102	55.9347	57.7221	55.0233	58.6335	-0.8284	0.592	-1.401	0.118		
5	56	56.4186	0.4894	55.3524	57.4849	54.5222	58.3151	-0.4186	0.528	-0.793	0.068		
6	56	55.8263	0.3354	55.0955	56.5572	54.0961	57.5566	0.1737	0.637	0.273	0.003		
7	55	55.2852	0.3695	54.4800	56.0903	53.5222	57.0481	-0.2852	0.618	-0.462	0.010		
8	55	55.1387	0.4487	54.1611	56.1163	53.2906	56.9867	-0.1387	0.563	-0.246	0.005		
9	55	54.8051	0.5010	53.7134	55.8967	52.8942	56.7159	0.1949	0.517	0.377	0.017		
10	55	55.1592	0.4663	54.1431	56.1752	53.2905	57.0278	-0.1592	0.548	-0.290	0.008		
11	54	54.2200	0.3731	53.4070	55.0329	52.4535	55.9864	-0.2200	0.616	-0.357	0.006		
12	54	53.9125	0.5763	52.6568	55.1681	51.9034	55.9215	0.0875	0.431	0.203	0.009		
13	54	53.1392	0.3714	52.3301	53.9484	51.3745	54.9040	0.8608	0.617	1.396	0.088		
14	53	53.1109	0.5091	52.0015	54.2202	51.1899	55.0319	-0.1109	0.509	-0.218	0.006		
15	53	53.8286	0.2863	53.2048	54.4524	52.1407	55.5164	-0.8286	0.660	-1.255	0.037		
16	53	52.8268	0.6023	51.5145	54.1391	50.7818	54.8718	0.1732	0.394	0.439	0.056		
17	52	51.7865	0.5118	50.6715	52.9015	49.8622	53.7108	0.2135	0.506	0.422	0.023		
18	52	51.4349	0.5011	50.3432	52.5266	49.5240	53.3458	0.5651	0.517	1.093	0.140		
19	51	51.4835	0.3672	50.6833	52.2836	49.7228	53.2441	-0.4835	0.619	-0.781	0.027		
20	50	50.4242	0.5004	49.3340	51.5144	48.5142	52.3342	-0.4242	0.517	-0.820	0.079		

Exercícios

- bi) Mostre que o valor do resíduo usual associado à primeira observação é $e_1=1.5105$
- bii) Mostre que o valor do resíduo standardizado associado à primeira observação é $R_1=2.577$

The SAS System

										Ei	Ri	Di	hii	Ti
Obs	Area	Gordurasubcut	Peso	Rendimento	Gordurarenalpel	Comprimento	LarguraAnca	Musculo	predicted	r	int_r	cooksd	leverage	ext_r
1	8.8	3.5	13.9	50	200	72	25	58	56.4895	1.51050	2.57743	0.42232	0.33713	3.69342
2	9.2	3.0	15.0	47	170	68	24	57	57.4744	-0.47437	-1.09829	0.26799	0.63994	-1.10874
3	8.6	3.5	13.4	48	180	73	23	57	56.4076	0.59236	0.91688	0.02536	0.19443	0.91031
4	8.7	4.0	14.2	48	150	74	25	56	56.8284	-0.82843	-1.40054	0.11790	0.32471	-1.46607
5	8.5	3.5	13.0	51	160	69	22	56	56.4186	-0.41865	-0.79311	0.06758	0.46223	-0.78007
6	8.2	4.0	14.8	49	190	70	21	56	55.8263	0.17368	0.27271	0.00258	0.21717	0.26191
7	8.0	4.5	12.8	46	210	71	27	55	55.2852	-0.28516	-0.46164	0.00953	0.26356	-0.44596
8	7.9	5.0	15.1	48	200	76	23	55	55.1387	-0.13867	-0.24637	0.00482	0.38853	-0.23648
9	7.6	4.5	13.6	47	190	65	20	55	54.8051	0.19494	0.37718	0.01671	0.48447	0.36328
10	7.5	5.0	14.1	50	210	66	28	55	55.1592	-0.15915	-0.29025	0.00762	0.41970	-0.27888
11	7.6	4.5	13.7	49	250	65	22	54	54.2200	-0.21996	-0.35732	0.00586	0.26867	-0.34395
12	7.4	4.0	12.2	48	280	74	21	54	53.9125	0.08754	0.20298	0.00920	0.64101	0.19468
13	7.3	6.0	14.7	49	230	68	20	54	53.1392	0.86076	1.39598	0.08837	0.26620	1.46038
14	7.0	5.5	13.1	52	280	62	26	53	53.1109	-0.11086	-0.21788	0.00594	0.50031	-0.20902
15	7.5	5.0	14.0	50	250	72	21	53	53.8286	-0.82858	-1.25462	0.03698	0.15820	-1.28869
16	6.8	6.0	12.9	43	300	75	29	53	52.8268	0.17321	0.43946	0.05638	0.70019	0.42418
17	6.5	6.5	14.2	53	310	71	23	52	51.7865	0.21348	0.42173	0.02272	0.50547	0.40680
18	6.8	7.0	12.8	45	260	65	22	52	51.4349	0.56507	1.09344	0.14049	0.48455	1.10329
19	7.0	6.5	13.5	47	290	68	20	51	51.4835	-0.48349	-0.78098	0.02683	0.26029	-0.76749
20	6.8	7.0	12.9	48	330	69	21	50	50.4242	-0.42420	-0.81980	0.07856	0.48323	-0.80785

Exercícios

biii) Mostre que o erro padrão do resíduo associado à primeira observação é 0.586

biv) Mostre que o valor da distância de Cook da primeira observação é $D_1=0.42232$

Algumas transformações de variáveis

Por vezes, é possível tornear violações às hipóteses de Normalidade dos erros aleatórios ou homogeneidade de variâncias através de transformações de variáveis. Por exemplo,

Se
$$var(\varepsilon_i) \propto E[Y_i]$$
 então $Y \longrightarrow \sqrt{Y}$
Se $var(\varepsilon_i) \propto (E[Y_i])^2$ então $Y \longrightarrow \ln Y$
Se $var(\varepsilon_i) \propto (E[Y_i])^4$ então $Y \longrightarrow 1/Y$

são propostas usuais para estabilizar as variâncias.

Os exemplos acima são casos particulares da família Box-Cox de transformações:

$$Y \longrightarrow \begin{cases} \frac{Y^{\lambda}-1}{\lambda} & , \ \lambda \neq 0 \\ \ln(Y) & , \ \lambda = 0 \end{cases}$$

Prevenções sobre transformações

Mas a utilização de transformações de variáveis, sobretudo quando afecta a variável resposta, deve ser feita com cautela.

- Uma transformação de variáveis muda também a relação de base entre as variáveis originais;
- Uma transformação que "corrija" um problema (e.g., variâncias heterogéneas) pode gerar outro (e.g., não-normalidade);
- Existe o perigo de usar transformações que resolvam o problema duma amostra específica, mas não tenham qualquer generalidade.

Advertências finais

- Podem surgir problemas associados à (quase) multicolinearidade das variáveis preditoras, ou seja, ao facto das colunas da matriz X serem (quase) linearmente dependentes:
 - podem existir problemas numéricos no cálculo de (X^tX)⁻¹, logo no ajustamento do modelo e na estimação dos parâmetros;
 - podem existir variâncias muito grandes de alguns $\hat{\beta}_i$ s, o que significa muita instabilidade na inferência.

Multicolinearidade reflecte redundância de informação nos preditores. É possível eliminá-la excluíndo da análise uma ou várias variáveis preditoras que sejam responsáveis pela (quase) dependência linear dos preditores.

Advertências finais (cont.)

2. Não se deve confundir a existência de uma relação linear entre preditores $X_1, X_2, ..., X_p$ e uma variável resposta Y, com uma relação de causa e efeito.

Pode existir uma relação de causa e efeito. Mas pode também verificar-se:

- Uma relação de variação conjunta, mas não de tipo causal (como por exemplo, em muitos conjuntos de dados morfométricos). Por vezes, preditores e variável resposta são todos efeito de causas comuns subjacentes.
- Uma relação espúria, de coincidência numérica.

Uma relação causal só pode ser afirmada com base em teoria própria do fenómeno sob estudo, e não com base na relação linear estabelecida estatisticamente.