
INSTITUTO SUPERIOR DE AGRONOMIA
MODELOS MATEMÁTICOS e APLICAÇÕES – 2020-21
Resoluções de exercícios de Modelo Linear

1 Regressão Linear

1. Escreva, numa sessão do R, o comando indicado no enunciado:

```
> Cereais <- read.csv("Cereais.csv")
```

Para ver o *conteúdo* do objecto `Cereais` acabado de criar, escrevemos o seu nome, como ilustrado de seguida (tendo sido omitidas várias linhas do conteúdo por razões de espaço):

```
> Cereais
  ano  area
1 1986 8789.69
2 1987 8972.11
3 1988 8388.94
4 1989 9075.35
5 1990 7573.48
(...)
24 2009 3398.99
25 2010 3041.18
26 2011 2830.96
```

NOTA: O comando `read.csv` parte do pressuposto que o ficheiro indicado contém colunas de dados - cada coluna correspondente a uma variável. O objecto `Cereais` criado no comando acima é uma *data frame*, que pode ser encarada como uma tabela de dados em que cada coluna corresponde a uma variável. As variáveis individuais da *data frame* podem ser acedidas através duma indexação análoga à utilizada para objectos de tipo matriz, referenciando o número da respectiva coluna:

```
> Cereais[,2]
 [1] 8789.69 8972.11 8388.94 9075.35 7573.48 8276.47 7684.20 7217.93 6773.54
[10] 6756.57 6528.18 6902.34 5065.38 5923.45 5779.21 4927.15 5149.21 4507.98
[19] 4636.46 3893.43 3731.92 3120.99 3653.74 3398.99 3041.18 2830.96
```

Alternativamente, as variáveis que compõem uma *data frame* podem ser acedidas através do nome da *data frame*, seguido dum cifrão e do nome da variável:

```
> Cereais$area
 [1] 8789.69 8972.11 8388.94 9075.35 7573.48 8276.47 7684.20 7217.93 6773.54
[10] 6756.57 6528.18 6902.34 5065.38 5923.45 5779.21 4927.15 5149.21 4507.98
[19] 4636.46 3893.43 3731.92 3120.99 3653.74 3398.99 3041.18 2830.96
```

(a) `> plot(Cereais)`

O gráfico obtido revela uma forte relação linear (decrecente) entre anos e superfície agrícola dedicada à produção de cereais.

Repare-se que o comando funciona correctamente nesta forma muito simples porque: (i) a *data frame* `Cereais` apenas tem duas variáveis; e (ii) a ordem dessas variáveis coincide com

a ordem desejada no gráfico: a primeira variável no eixo horizontal e a segunda no eixo vertical.

Existe uma forma mais geral do comando que também poderia ser usada neste caso: `plot(x,y)`, onde `x` e `y` indicam os nomes das variáveis que desejamos ocupar, respectivamente o eixo horizontal e o eixo vertical. No nosso exemplo, poderíamos escrever:

```
> plot(Cereais$ano, Cereais$area)
```

- (b) O gráfico obtido na alínea anterior apresenta uma tendência linear decrescente, pelo que o coeficiente de correlação será negativo. A tendência linear é bastante acentuada, pelo que é de supor que o coeficiente de correlação seja próximo de -1 .

O comando `cor` do R calcula coeficientes de correlação. Se os seus argumentos forem dois vectores (necessariamente de igual dimensão), é devolvido o coeficiente de correlação. Se o seu argumento for uma *data frame*, é devolvida uma matriz de correlações entre todos os pares de variáveis da *data frame*. No nosso caso, esta segunda alternativa produz:

```
> cor(Cereais)
           ano      area
ano  1.0000000 -0.9826927
area -0.9826927  1.0000000
```

O coeficiente de correlação entre `ano` e `area` é, como previsto, muito próximo de -1 , confirmando a existência duma forte relação linear decrescente entre anos e superfície agrícola para a produção de cereais em Portugal, nos anos indicados.

- (c) Os parâmetros da recta podem ser calculados, quer a partir da sua definição, quer utilizando o comando do R que ajusta uma regressão linear: o comando `lm` (as iniciais, pela ordem em inglês, de *modelo linear*). Sabemos que: $b_1 = \frac{cov_{xy}}{s_x^2}$ e $b_0 = \bar{y} - b_1 \bar{x}$. Utilizando o R, é possível calcular os indicadores estatísticos nas definições:

```
> cov(Cereais$ano, Cereais$area)
[1] -15137.48
> var(Cereais$ano)
[1] 58.5
> -15137.48/58.5
[1] -258.7603
> mean(Cereais$area)
[1] 5869.187
> mean(Cereais$ano)
[1] 1998.5
> 5869.187 - (-258.7603)*1998.5
[1] 523001.6
```

O comando `lm` devolve directamente os parâmetros da recta de regressão:

```
> lm(area ~ ano, data=Cereais)
Call:
lm(formula = area ~ ano, data = Cereais)
Coefficients:
(Intercept)      ano
 523001.7      -258.8
```

NOTA: Na fórmula $y \sim x$, a variável do lado esquerdo do til é a variável resposta, e a do lado direito é a variável preditora. O argumento `data` permite indicar o objecto onde se encontram as variáveis cujos nomes são referidos na fórmula.

O resultado do ajustamento pode ser guardado como um novo objecto, que poderá ser invocado sempre que se deseje trabalhar com a regressão agora ajustada:

```
> Cereais.lm <- lm(area ~ ano, data=Cereais)
```

Interpretação dos coeficientes:

- Declive: $b_1 = -258.8 \text{ km}^2/\text{ano}$ indica que, em cada ano que passa, a superfície agrícola dedicada à produção de cereais diminui, em média, $258,8 \text{ km}^2$. Em geral (e como se pode comprovar analisando a fórmula para o declive da recta de regressão), as unidades de b_1 são as unidades da variável resposta y a dividir pelas unidades da variável preditora x . Fala-se em “variação média” porque a recta apenas descreve a tendência de fundo, na relação entre x e y .
 - Ordenada na origem: $b_0 = 523001.7 \text{ km}^2$. Em geral, as unidades de b_0 são as unidades da variável resposta y . A interpretação deste valor é bizarra: a superfície agrícola utilizada na produção de cereais no ano $x = 0$, seria cerca de 5 vezes superior à área total do país, situação claramente impossível. A impossibilidade ilustra a ideia geral de que, *na ausência de mais informação, a validade dum relação linear não poder ser extrapolada para fora da gama de valores de x observada* (neste caso, os anos 1986-2011).
- (d) Sabe-se que, numa regressão linear simples entre variáveis x e y , o coeficiente de determinação é o quadrado do coeficiente de correlação entre as variáveis, ou seja: $R^2 = r_{xy}^2$. O valor do coeficiente de correlação entre x e y pode ser obtido através do comando `cor`:

```
> cor(Cereais$ano, Cereais$area)
[1] -0.9826927
> cor(Cereais$ano, Cereais$area)^2
[1] 0.9656849
```

No nosso caso $R^2 = 0.9656849$, ou seja, cerca de 96,6% da variabilidade total observada para a variável resposta y é explicada pela regressão.

O comando `summary`, aplicando ao resultado da regressão ajustada, produz vários resultados de interesse relativos à regressão. O coeficiente de determinação pedido nesta alínea é indicado na penúltima linha da listagem produzida:

```
> summary(Cereais.lm)
(...)
Multiple R-squared: 0.9657
(...)
```

- (e) O comando `abline(Cereais.lm)` traça a recta pedida em cima do gráfico anteriormente criado pelo comando `plot`. Confirma-se o bom ajustamento da recta à nuvem de pontos, já indiciado pelo valor muito elevado do R^2 .

Nota: Em geral, o comando `abline(a,b)` traça, num gráfico já criado, a recta de equação $y = a + bx$. No caso do `input` ser o ajustamento dum regressão linear simples (obtido através do comando `lm` e que devolve o par de coeficientes b_0 e b_1), o resultado é o gráfico da recta $y = b_0 + b_1 x$.

- (f) Sabemos que $SQT = (n-1) s_y^2$, pelo que podemos calcular este valor através do comando:

```
> (length(Cereais$area)-1)*var(Cereais$area)
[1] 101404176
```

- (g) Sabemos que $R^2 = \frac{SQR}{SQT}$, pelo que $SQR = R^2 \times SQT$:

```
> 0.9656849*101404176
[1] 97924482
```

Alternativamente, e uma vez que $SQR = (n-1) s_{\hat{y}}^2$, pode-se usar o comando `fitted` para obter os valores ajustados de y (\hat{y}_i) e seguidamente obter o valor de SQR :

```
> fitted(Cereais.lm)
      1      2      3      4      5      6      7      8
9103.691 8844.930 8586.170 8327.410 8068.649 7809.889 7551.129 7292.368
      9     10     11     12     13     14     15     16
7033.608 6774.848 6516.087 6257.327 5998.567 5739.806 5481.046 5222.286
(...)
> (length(Cereais$area)-1)*var(fitted(Cereais.lm))
[1] 97924480
```

NOTA: A discrepância nos dois valores obtidos para SQR deve-se a erros de arredondamento.

- (h) O comando `residuals` devolve os resíduos dum modelo ajustado. Logo,

```
> residuals(Cereais.lm)
      1      2      3      4      5      6      7
-314.00068 127.17965 -197.23002 747.94031 -495.16936 466.58097 133.07131
      8      9     10     11     12     13     14
-74.43836 -260.06803 -18.27770 12.09263 645.01296 -933.18670 183.64363
(...)
> sum(residuals(Cereais.lm)^2)
[1] 3479697
```

É fácil de verificar que se tem $SQR + SQRE = SQT$:

```
> 97924480+3479697
[1] 101404177
```

- (i) Com o auxílio do R, podemos efectuar o novo ajustamento. No caso de se efectuar uma transformação duma variável, esta deve ser efectuada, na fórmula do comando `lm`, com a protecção `I()`, como indicado no comando seguinte:

```
> lm(I(area*100) ~ ano, data=Cereais)
Call:
lm(formula = I(area * 100) ~ ano, data = Cereais)
Coefficients:
(Intercept)      ano
 52300171      -25876
```

Comparando estes valores dos parâmetros ajustados com os que haviam sido obtidos inicialmente, verifica-se que ambos os parâmetros ajustados foram multiplicados por 100. Não se trata duma coincidência, como se pode verificar analisando o efeito da transformação $y \rightarrow y^* = cy$ (para qualquer constante c) nas fórmulas dos parâmetros da recta ajustada. Indicando por b_1 e b_0 os parâmetros na recta original e por b_1^* e b_0^* os novos parâmetros, obtidos com a transformação indicada, temos (recordando que $cov(x, cy) = c cov(x, y)$):

$$b_1^* = \frac{cov_x y^*}{s_x^2} = \frac{cov(x, cy)}{s_x^2} = c \frac{cov(x, y)}{s_x^2} = c b_1 ;$$

e (tendo em conta o efeito de constantes multiplicativas sobre a média, ou seja, $\overline{y^*} = c\overline{y}$):

$$b_0^* = \overline{y^*} - b_1^* \overline{x} = c\overline{y} - c b_1 \overline{x} = c(\overline{y} - b_1 \overline{x}) = c b_0 .$$

Assim, multiplicar a variável resposta por uma constante c tem por efeito multiplicar os dois parâmetros da recta ajustada por essa mesma constante c . No entanto, o coeficiente de determinação permanece inalterado. Esse facto, que resulta da invariância do valor absoluto do coeficiente de correlação a qualquer transformação linear de uma, ou ambas as variáveis, pode ser confirmado através do R:

```
> summary(lm(I(area*100) ~ ano, data=Cereais))
(...)
Multiple R-squared: 0.9657
(...)
```

- (j) Nesta alínea é pedida uma translação da variável preditora, da forma $x \rightarrow x^* = x + a$, com $a = -1985$. Neste caso, e comparando com o ajustamento inicial, verifica-se que o declive da recta de regressão não se altera, mas a sua ordenada na origem sim:

```
> lm(area ~ I(ano-1985), data=Cereais)
Call:
lm(formula = area ~ I(ano - 1985), data = Cereais)
Coefficients:
(Intercept)  I(ano - 1985)
    9362.5         -258.8
```

Inspeccionando o efeito duma translação na variável preditora sobre o declive da recta ajustada, temos (tendo em conta que constantes aditivas não alteram, nem a variância, nem a covariância):

$$b_1^* = \frac{\text{cov}_y x^*}{s_{x^*}^2} = \frac{\text{cov}(x, y)}{s_x^2} = b_1 .$$

Já no que respeita à ordenada na origem, e tendo em conta a forma como os valores médios são afectados por constantes aditivas, tem-se:

$$b_0^* = \bar{y} - b_1^* \bar{x}^* = \bar{y} - b_1 (\bar{x} + a) = (\bar{y} - b_1 \bar{x}) - b_1 a = b_0 - a b_1 .$$

Assim, a nova ordenada na origem vem $b_0^* = 523001.6 - (-1985) * (-258.7603) = 9362.405$. Tal como na alínea anterior, a transformação da variável preditora é linear, pelo que o coeficiente de determinação não se altera: $R^2 = 0.9657$.

2. (a) Seguindo as instruções do enunciado, cria-se a *data frame* **azeite**, com os dados:

```
> library(xlsx)
> azeite <- read.xlsx("Azeite.xls", sheetIndex=1, header=TRUE)
```

Notas:

- i. Caso o ficheiro **Azeite.xls** esteja numa directoria diferente da directoria de trabalho do R, o nome do ficheiro deverá incluir a sequência de pastas e subpastas que devem ser percorridas para chegar até ao ficheiro.
- ii. O argumento **header** tem valor lógico que indica se a primeira linha do ficheiro a ser lido contém, ou não, os nomes das variáveis. Por omissão o argumento tem o valor lógico **FALSE**, que considera que na primeira linha do ficheiro já há valores numéricos.

O resultado do comando pode ser visto escrevendo o nome do objecto agora lido:

```
> azeite
  Ano Azeitona Azeite
1 1995   311257 477728
2 1996   275143 452038
3 1997   309090 423584
4 1998   225616 360948
5 1999   320865 512264
6 2000   167161 249433
7 2001   218522 349502
```

```

8 2002 211574 310474
9 2003 232947 364976
10 2004 300699 500658
11 2005 203909 318174
12 2006 362301 518466
13 2007 203968 352574
14 2008 336479 587422
15 2009 414687 681850
16 2010 435009 686832

```

- (b) Quando aplicado a uma *data frame*, o comando `plot` produz uma “matriz de gráficos” de cada possível par de variáveis (confirme!). Neste caso, não é pedido qualquer gráfico envolvendo a primeira variável da *data frame* (Ano). Existem várias maneiras alternativas de pedir apenas o gráfico das segunda e terceira variáveis, uma das quais envolve o conceito de *indexação negativa*, que tanto pode ser utilizado em *data frames* como em matrizes: índices negativos representam linhas ou colunas a serem *omitidas*. Assim, qualquer dos seguintes comandos (alternativos) produz o gráfico pedido no enunciado:

```

> plot(azeite[,-1])
> plot(azeite[,c(2,3)])
> plot(azeite$Azeitona, azeite$Azeite)

```

- (c) O comando `cor` do R calcula a matriz dos coeficientes de correlação entre cada par de variáveis da *data frame*.

```

> cor(azeite)
      Ano Azeitona  Azeite
Ano    1.0000000 0.3999257 0.4715217
Azeitona 0.3999257 1.0000000 0.9722528
Azeite  0.4715217 0.9722528 1.0000000

```

O valor da correlação pedido é $r_{xy} = 0.9722528$, um valor positivo muito elevado, que indica uma relação linear crescente muito forte, entre produção de azeitona e produção de azeite.

- (d) Utilizando o comando `lm` do R, tem-se:

```

> lm(Azeite ~ Azeitona, data=azeite)
Call: lm(formula = Azeite ~ Azeitona, data = azeite)
Coefficients:
(Intercept)      Azeitona
   -5151.793         1.596

```

A cada tonelada adicional de produção de azeitona oleificada corresponde um aumento médio de 1.596 *hl* de produção de azeite. O valor da ordenada na origem indica uma situação impossível: na ausência de produção de azeitona, a produção média de azeite seria negativa ($b_0 = -5151.793$ *hl*). O modelo não deve ser usado para produções de azeitona próximas de zero. Como sempre, deve ser usado com cautela fora da gama de valores observados de x .

- (e) “*Precisão da recta*” é uma designação alternativa para o coeficiente de determinação R^2 . Sabe-se que numa regressão linear simples, $R^2 = r_{xy}^2$. Logo, $R^2 = 0.9722528^2 = 0.9452755$. Cerca de 94.5% da variabilidade na produção de azeite é explicável pela regressão linear simples sobre a produção de azeitona.

$$3. \quad (a) \quad \sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} = n\bar{x} - n\bar{x} = 0.$$

- (b) Por definição, $(n-1)cov_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$. Distribuindo o primeiro factor de cada parcela pelas parcelas do segundo factor e utilizando o resultado da alínea anterior, temos:

$$(n-1)cov_{xy} = \sum_{i=1}^n (x_i - \bar{x})y_i - \sum_{i=1}^n (x_i - \bar{x})\bar{y} = \sum_{i=1}^n (x_i - \bar{x})y_i - \bar{y} \underbrace{\sum_{i=1}^n (x_i - \bar{x})}_{=0} = \sum_{i=1}^n (x_i - \bar{x})y_i$$

Trocando o papel das variáveis x e y , mostra-se que $(n-1)cov_{xy} = \sum_{i=1}^n x_i(y_i - \bar{y})$.

A partir da expressão final acima, e distribuindo x_i pela diferença de cada parcela, vem:

$$(n-1)cov_{xy} = \sum_{i=1}^n (x_i - \bar{x})y_i = \sum_{i=1}^n x_i y_i - \sum_{i=1}^n \bar{x} y_i = \sum_{i=1}^n x_i y_i - \bar{x} \underbrace{\sum_{i=1}^n y_i}_{=n\bar{y}} = \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}.$$

4. (a) Tendo em conta que os valores ajustados de y são dados por $\hat{y}_i = b_0 + b_1 x_i$, tem-se que a média dos valores ajustados é dada por:

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

Mas a ordenada de origem duma recta de regressão é dada por $b_0 = \bar{y} - b_1 \bar{x}$, pelo que a última expressão equivale à média \bar{y} dos valores observados de y .

- (b) Tem-se, por definição, que $e_i = y_i - \hat{y}_i$. Logo (e tendo em conta a alínea anterior),

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

- (c) Pela definição de coeficiente de correlação entre x e y , tem-se:

$$r_{xy} = \frac{cov_{xy}}{s_x \cdot s_y} = \frac{cov_{xy}}{s_x^2} \cdot \frac{s_x}{s_y} = b_1 \cdot \frac{s_x}{s_y}$$

- (d) Por definição, $R^2 = \frac{SQR}{SQT}$. Sabemos que $SQT = (n-1)s_y^2$. Verifiquemos que $SQR = b_1^2(n-1)s_x^2$. De facto, recordando a definição dos valores ajustados de y e a expressão da ordenada na origem da recta de regressão, b_0 , temos que $\hat{y}_i = b_0 + b_1 x_i = \bar{y} + b_1(x_i - \bar{x})$. Logo,

$$SQR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n [b_1(x_i - \bar{x})]^2 = b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 = b_1^2(n-1)s_x^2.$$

Tendo também em conta o resultado da alínea anterior, tem-se $R^2 = \frac{b_1^2 s_x^2}{s_y^2} = (r_{xy})^2$.

- (e) Os valores ajustados \hat{y}_i são dados por uma mesma transformação linear (afim) dos valores do preditor: $\hat{y}_i = b_0 + b_1 x_i$. São conhecidas as propriedades destas transformações sobre a covariância e a variância. Assim,

$$r_{y\hat{y}}^2 = \left(\frac{cov_{y\hat{y}}}{s_y s_{\hat{y}}} \right)^2 = \frac{cov_{y, b_0 + b_1 x}^2}{s_y^2 \cdot s_{b_0 + b_1 x}^2} = \frac{(b_1 cov_{y,x})^2}{s_y^2 \cdot b_1^2 s_x^2} = \frac{b_1^2 \cdot cov_{xy}^2}{b_1^2 s_x^2 \cdot s_y^2} = r_{xy}^2 = R^2.$$

Assim, o coeficiente de determinação duma regressão linear simples é também o quadrado do coeficiente de correlação linear entre valores observados e ajustados de y . Esta propriedade estende-se às regressões lineares múltiplas, embora seja necessário adaptar a justificação.

5. Neste exercício é pedido para discutir a recta forçada à origem relacionando, nos dados *iris*, as variáveis relativas às pétalas. No Exercício 14 é considerada a recta usual destes mesmos dados.

- (a) Qualquer que venha a ser o valor do parâmetro b , os valores ajustados serão $\hat{y}_i = b x_i$ e os resíduos $e_i = y_i - \hat{y}_i$. Logo, minimizar a Soma de Quadrados dos Resíduos corresponde a minimizar a função $SQRE(b) = \sum_{i=1}^n (y_i - b x_i)^2$. Trata-se duma função duma única variável (b), e condição necessária para um seu ponto mínimo será anular a derivada. Tem-se:

$$\frac{dSQRE}{db} = 2 \sum_{i=1}^n (y_i - b x_i)(-x_i) = 0 \Leftrightarrow -\sum_{i=1}^n y_i x_i + b \sum_{i=1}^n x_i^2 = 0 \Leftrightarrow b = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}.$$

Assinale-se que na recta de regressão usual (não forçada à origem), o declive é dado por $b_1 = \frac{cov_{xy}}{s_x^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$. A expressão agora obtida para o declive da recta forçada à origem substitui os momentos centrados pelos correspondentes momentos não centrados.

- (b) A regressão forçada à origem pedida no enunciado é:

```
> irisFT0 <- lm(Petal.Width ~ -1 + Petal.Length , data=iris)
> coef(irisFT0)
Petal.Length
0.3365109
```

A fórmula para b da alínea anterior confirma o valor obtido:

```
> sum(iris$Petal.Length*iris$Petal.Width)/sum(iris$Petal.Length^2)
[1] 0.3365109
```

Nesta regressão:

- i. A soma dos resíduos não é nula:

```
> sum(residuals(irisFT0))
[1] -9.79118
```

De facto, $\sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - b x_i) = \sum_{i=1}^n y_i - b \sum_{i=1}^n x_i = n\bar{y} - b n\bar{x}$. Substituindo a expressão para b obtida na alínea anterior, obtém-se uma expressão que não é, em geral, zero.

- ii. Se já não é verdade que \bar{e} seja zero, a Soma de Quadrados dos Resíduos, $\sum_{i=1}^n e_i^2$ já não é igual ao numerador da variância dos resíduos, $(n-1) s_e^2 = \sum_{i=1}^n (e_i - \bar{e})^2$. Confirme-se para os dados em questão:

```
> sum(residuals(irisFT0)^2)
[1] 9.865034
> 149*var(residuals(irisFT0))
[1] 9.225919
```

- iii. As expressões usuais das Somas de Quadrados são:

$$SQT = \sum_{i=1}^n (y_i - \bar{y})^2 \quad ; \quad SQR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad ; \quad SQRE = \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

No exemplo em apreço, tem-se:


```

> sqt <- 149*var(iris$Petal.Width)
> sqr <- sum((fitted(irisFT0)-mean(iris$Petal.Width))^2)
> sqre <- sum((iris$Petal.Width - fitted(irisFT0))^2)
> sqt
[1] 86.56993
> sqr+sqre
[1] 63.08416

```

A principal conclusão a retirar é que as quantidades habitualmente utilizadas na regressão linear usual não podem ser mecanicamente transportadas para o estudo duma recta forçada à origem. A igualdade $SQT = SQR + SQRE$ estava na base da própria definição do coeficiente de determinação que, na regressão usual, tanto podia ser definido como $\frac{SQR}{SQT}$ ou pela expressão equivalente $1 - \frac{SQRE}{SQT}$. *Numa recta forçada à origem estas duas expressões não são equivalentes, nem é legítimo afirmar que o quociente $\frac{SQR}{SQT}$ represente a proporção da variabilidade total explicada pela regressão*, uma vez que já não é verdade que SQR seja o numerador da variância dos \hat{y}_i (cuja média, como se viu, já não é \bar{y}). O próprio conceito de R^2 tem de ser adaptado neste contexto. Tudo isto resulta da exclusão, que se poderia erradamente pensar ser trivial, da constante aditiva na equação da recta.

6. Os dados referidos no enunciado são obtidos como se indica a seguir:

```

> library(MASS)
> Animals

```

	body	brain
Mountain beaver	1.350	8.1
Cow	465.000	423.0
Grey wolf	36.330	119.5
Goat	27.660	115.0
Guinea pig	1.040	5.5
Dipliodocus	11700.000	50.0
Asian elephant	2547.000	4603.0
Donkey	187.100	419.0
Horse	521.000	655.0
Potar monkey	10.000	115.0
Cat	3.300	25.6
Giraffe	529.000	680.0
Gorilla	207.000	406.0
Human	62.000	1320.0
African elephant	6654.000	5712.0
Triceratops	9400.000	70.0
Rhesus monkey	6.800	179.0
Kangaroo	35.000	56.0
Golden hamster	0.120	1.0
Mouse	0.023	0.4
Rabbit	2.500	12.1
Sheep	55.500	175.0
Jaguar	100.000	157.0
Chimpanzee	52.160	440.0
Rat	0.280	1.9
Brachiosaurus	87000.000	154.5
Mole	0.122	3.0
Pig	192.000	180.0

(a) A nuvem de pontos pedida pode ser obtida através do comando `plot(Animals)`. Quanto ao coeficiente de correlação, tem-se:

```
> cor(Animals)
              body      brain
body  1.000000000 -0.005341163
brain -0.005341163  1.000000000
```

O valor quase nulo do coeficiente de correlação indica ausência de relacionamento linear entre os pesos do corpo e do cérebro, facto que se confirma visualmente no gráfico.

- (b) É pedida a nuvem de pontos das transformações logarítmicas das duas variáveis da *data frame* `Animals`, que pode ser obtida duma das seguintes formas:

```
> plot(log(Animals))
ou, alternativamente,
> plot(log(brain) ~ log(body), data=Animals)
```

NOTA: Os logaritmos aqui referidos são os logaritmos naturais, \ln . Por omissão, o comando `log` do R calcula logaritmos naturais.

Os coeficientes de correlação e de determinação entre log-pesos do corpo e log-pesos do cérebro podem ser calculados, com o auxílio do R, da seguinte forma:

```
> cor(log(Animals$body), log(Animals$brain))    <-- coeficiente de correlação
[1] 0.7794935
> cor(log(Animals$body), log(Animals$brain))^2  <-- coeficiente de determinação
[1] 0.6076101
```

Como $R^2 = 0.6076$, a regressão linear entre log-peso do cérebro e log-peso do corpo explica menos de 61% da variabilidade total *dos log-pesos* do cérebro observados. Este valor, aparentemente contraditório com a relativamente forte relação linear para a maioria das espécies, é reflexo da presença nos dados de três espécies claramente atípicas face às restantes.

- (c) Como se viu nas aulas, uma relação linear entre $\ln(y)$ e $\ln(x)$ corresponde a uma relação potência (alométrica) entre as variáveis originais: $y = cx^d$. Neste caso, tem-se uma relação de tipo alométrico entre pesos duma parte do organismo (cérebro) e do todo (corpo). O gráfico indica que é aceitável admitir uma relação potência entre o peso do cérebro e o peso do corpo, nas espécies animais consideradas.

- (d) Os comandos pedidos são:

```
> Animals.loglm <- lm(log(brain) ~ log(body), data=Animals)
> Animals.loglm
Call: lm(formula = log(brain) ~ log(body), data = Animals)
Coefficients:
(Intercept)    log(body)
      2.555         0.496
> abline(Animals.loglm)
```

- (e) O declive $b_1^* = 0.49599$ da recta ajustada tem duas leituras possíveis. Na relação entre as variáveis logaritmizadas tem a habitual leitura de qualquer declive duma recta de regressão: o log-peso do cérebro aumenta em média 0.49599 log-gramas, por cada aumento de 1 log-kg no peso do corpo. Mais compreensível é a interpretação na relação potência entre as variáveis originais. Como se viu nas aulas teóricas, a relação original entre y e x é da forma $y = cx^d$ com $d = b_1^* = 0.49599$ e $b_0^* = \ln(c) = 2.555 \Leftrightarrow c = e^{2.555} = 12.871$. No nosso caso, a tendência de fundo na relação entre peso do corpo (x) e peso do cérebro (y) é $y = 12.871 x^{0.49599}$. O valor de d muito próximo de 0.5 permite simplificar a relação dizendo que o ajustamento indica que o peso do cérebro é aproximadamente proporcional à *raíz quadrada* do peso do corpo.

- (f) O comando `identify(log(Animals))` permite, com o auxílio do rato, identificar pontos seleccionados pelo utilizador (para terminar, clique no botão direito do rato).

NOTA: É necessário explicitar as coordenadas dos pontos no gráfico que se vai aceder com o comando. No nosso caso, isso significa explicitar as coordenadas dos dados logaritizados: `log(Animals)`.

O enunciado pede para identificar os pontos que se destacam da relação linear: os pontos 6, 16 e 26. Seleccionando as linhas com esses números podemos identificar as espécies em questão, e verificar que se trata de espécies de dinossáurios:

```
> Animals[c(6,16,26),]
      body brain
Dipliodocus 11700 50.0
Triceratops  9400 70.0
Brachiosaurus 87000 154.5
```

- (g) Utilizando a indexação negativa para eliminar as três espécies de dinossáurios pode proceder-se ao reajustamento da regressão, modificando o argumento `data` do comando `lm`. O ajustamento sem as espécies extintas produz os seguintes parâmetros da recta:

```
> Animals.loglm.sub <- lm(log(brain) ~ log(body), data=Animals[-c(6,16,26),])
> Animals.loglm.sub
Call: lm(formula = log(brain) ~ log(body), data = Animals[-c(6,16,26),])
Coefficients:
(Intercept)      log(body)
      2.1504         0.7523
```

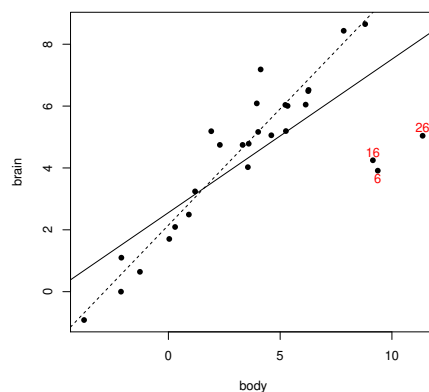
Note-se como o declive da recta se altera visivelmente, crescendo para mais de 0.75. Podemos analisar o efeito sobre o coeficiente de determinação aplicando o comando `summary` à regressão agora ajustada:

```
> summary(Animals.loglm.sub)$r.sq
[1] 0.9216991
```

Com a exclusão das espécies extintas, a recta de regressão passa a explicar mais de 92% da variabilidade total nos restantes log-pesos do cérebro, a partir dos log-pesos do corpo.

Pode juntar-se a nova recta ao gráfico obtido antes, através do comando `abline` (com um argumento pedindo uma recta a tracejado, para melhor a distinguir da recta original):

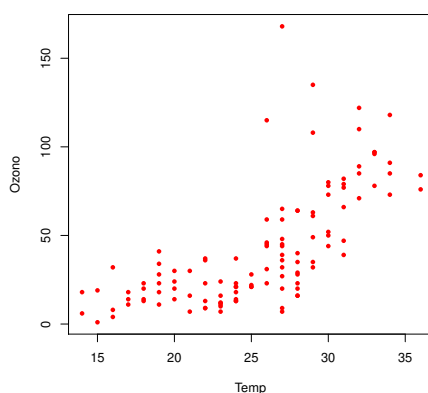
```
> abline(Animals.loglm.sub, lty="dashed")
```



A exclusão das três espécies de dinossáurios (as observações atípicas) permitiu que a recta ajustada acompanhe melhor a relação linear existente entre a generalidade das espécies do conjunto de dados. Este exemplo ilustra que *as rectas de regressão são sensíveis à presença de observações atípicas*.

- (h) O significado biológico dos parâmetros da recta é semelhante ao que foi visto na alínea 6e). Assim, na relação alométrica entre peso do cérebro e peso do corpo (variáveis não transformadas), o expoente será aproximadamente 0.75, ou seja, o peso do cérebro é aproximadamente proporcional à potência 3/4 do peso do corpo.
7. (a) O comando `plot(ozono)` produz o gráfico pedido. Um gráfico com alguns embelezamentos adicionais é produzido pelo comando:

```
> plot(ozono, col="red", pch=16, cex=0.8)
```



- (b) A linearização duma relação exponencial faz-se logaritmando:

$$y = ae^{bx} \Leftrightarrow \ln(y) = \ln(a) + bx,$$

que é uma relação linear entre x e $y^* = \ln(y)$.

- i. O gráfico de log-Ozono contra Temp pode ser construído pelo comando:

```
> plot(ozono$Temp, log(ozono$Ozono))
```

Uma tendência linear mais ou menos forte neste gráfico indica que a relação exponencial entre as variáveis originais é adequada. Neste caso, o gráfico corresponde a um coeficiente de correlação entre Temp e log-Ozono de 0.73.

- ii. O ajustamento pedido faz-se da seguinte forma:

```
> lm(log(Ozono) ~ Temp, data=ozono)
```

```
Call: lm(formula = log(Ozono) ~ Temp, data = ozono)
```

```
Coefficients:
```

```
(Intercept)      Temp
    0.3558         0.1203
```

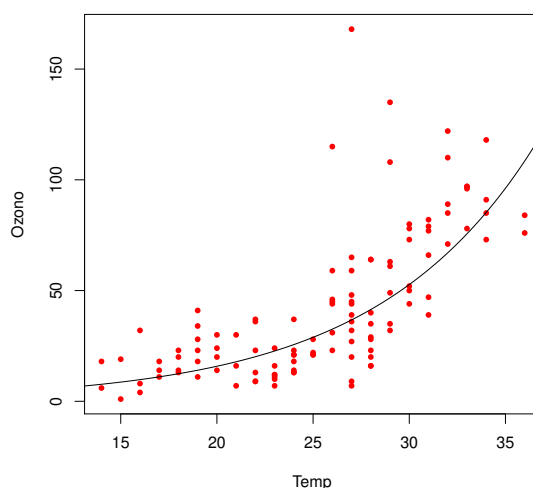
O coeficiente de determinação é de cerca de $R^2 = 0.73^2 = 0.53$ (aplicando o comando `summary` ao modelo agora ajustado verifica-se ser $R^2 = 0.5372$), o que significa que a regressão explica pouco mais de 53% da variabilidade dos log-teores de ozono.

- iii. O declive estimado da recta $b_1 = 0.1203$ é o coeficiente do expoente (b), na relação exponencial original. A ordenada na origem da recta ajustada, $b_0 = 0.3558$, corresponde

à estimativa de $\ln(a)$, pelo que a constante multiplicativa a da relação exponencial original é: $a = e^{0.3558} = 1.4273$.

- iv. a recta relaciona log-ozono com temperatura. Logo, o valor *de log-ozono* previsto pela recta, para um dia com temperatura máxima de 25° é dado por: $\widehat{y^*} = \widehat{\ln(y)} = 0.3558 + 0.1203 \times 25 = 3.3633$. E o teor estimado *de ozono* (em ppm) é: $e^{3.3633} = 28.8843$.
- (c) Eis um comando que ajusta a curva exponencial à nuvem de pontos de ozono vs. temperaturas (admitindo que este gráfico ainda está activo):

```
> curve(1.4273*exp(0.1203*x), from=10, to=40, add=TRUE)
```



8. (a) Com as restrições indicadas no enunciado, y não se anula e pode tomar-se o recíproco de y :

$$\frac{1}{y} = \frac{b+x}{ax} = \frac{b}{a} \cdot \frac{1}{x} + \frac{1}{a} \Leftrightarrow y^* = b_0^* + b_1^* x^*,$$

com $y^* = \frac{1}{y}$, $x^* = \frac{1}{x}$, $b_0^* = \frac{1}{a}$ e $b_1^* = \frac{b}{a}$. Assim, uma *relação linear entre os recíprocos de y e de x corresponde a uma relação de Michaelis-Menten entre y e x* .

- (b) Tendo em conta os nomes indicados no enunciado, e o facto de os dados do enunciado corresponderem *apenas às 12 primeiras linhas da data frame* (associadas ao valor `treated` da terceira coluna, de nome `state`), o modelo linearizado ajusta-se através do comando:

```
> lm(I(1/rate) ~ I(1/conc), data=Puromycin[Puromycin$state=="treated",])
```

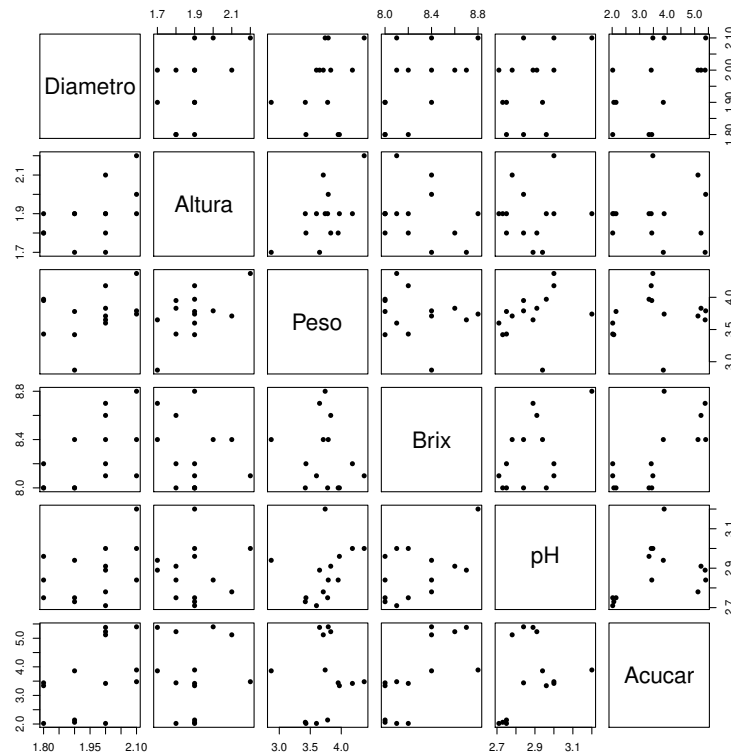
sendo os resultados obtidos os seguintes:

```
Coefficients:
(Intercept)    I(1/conc)
  0.0051072    0.0002472
```

- (c) Tendo em conta as relações vistas na alínea anterior, $b_0^* = \frac{1}{a} = 0.0051072$, tem-se $a = 195.802$. Por outro lado, $b_1^* = \frac{b}{a} = 0.0002472$, logo $b = 0.0002472 \times 195.802 = 0.04840225$. Assim, o modelo de Michaelis-Menten ajustado é: $y = \frac{195.802x}{0.04840225+x}$. Repare-se que o limite de y quando x tende para $+\infty$ é 195.802, que é assim a estimativa da assíntota superior da relação de Michaelis-Menten. O gráfico da relação original sugere que se pode tratar duma

subestimação do verdadeiro valor desta assintota horizontal. Este exemplo ilustra que pode haver inconvenientes associados à utilização de transformações linearizantes, como indicado nos acetatos das aulas teóricas.

9. (a) As nuvens de pontos e a matriz de correlações entre as variáveis da *data frame* `brix` são:



```
> round(cor(brix),d=3)
      Diametro  Altura   Peso   Brix   pH  Acucar
Diametro  1.000  0.488  0.302  0.557  0.411  0.492
Altura    0.488  1.000  0.587 -0.247  0.048  0.023
Peso      0.302  0.587  1.000 -0.198  0.308  0.118
Brix      0.557 -0.247 -0.198  1.000  0.509  0.714
pH        0.411  0.048  0.308  0.509  1.000  0.353
Acucar    0.492  0.023  0.118  0.714  0.353  1.000
```

É evidente nos gráficos a pequena dimensão do conjunto de dados. Conclui-se que não há relações lineares particularmente evidentes, facto confirmado pela matriz de correlações, onde a maior correlação é 0.714 (entre `Brix` e `Acucar`).

- (b) A equação de base (usando os nomes das variáveis como constam da *data frame*) é:

$$Brix_i = \beta_0 + \beta_1 Diametro_i + \beta_2 Altura_i + \beta_3 Peso_i + \beta_4 pH_i + \beta_5 Acucar_i + \epsilon_i ,$$

havendo nesta equação seis parâmetros (os cinco coeficientes β_j , $j = 1, 2, 3, 4, 5$, das variáveis predictoras e ainda a constante aditiva β_0).

- (c) Recorrendo ao comando `lm` do R, tem-se:

```

> brix.lm <- lm(Brix ~ . , data=brix)
> brix.lm
Call:
lm(formula = Brix ~ Diametro + Altura + Peso + pH + Acucar, data = brix)
Coefficients:
(Intercept)    Diametro      Altura        Peso          pH          Acucar
  6.08878      1.27093     -0.70967     -0.20453     0.51557     0.08971

```

(d) Tem-se:

```

> X <- model.matrix(brix.lm)
> X
  (Intercept) Diametro  Altura  Peso   pH  Acucar
1            1      2.0    2.1  3.71  2.78  5.12
2            1      2.1    2.0  3.79  2.84  5.40
3            1      2.0    1.7  3.65  2.89  5.38
4            1      2.0    1.8  3.83  2.91  5.23
5            1      1.8    1.8  3.95  2.84  3.44
6            1      2.0    1.9  4.18  3.00  3.42
7            1      2.1    2.2  4.37  3.00  3.48
8            1      1.8    1.9  3.97  2.96  3.34
9            1      1.8    1.8  3.43  2.75  2.02
10           1      1.9    1.9  3.78  2.75  2.14
11           1      1.9    1.9  3.42  2.73  2.06
12           1      2.0    1.9  3.60  2.71  2.02
13           1      1.9    1.7  2.87  2.94  3.86
14           1      2.1    1.9  3.74  3.20  3.89

```

A matriz do modelo é a matriz de dimensões $n \times (p+1)$, cuja primeira coluna é uma coluna de n uns e cujas p colunas seguintes são as colunas dadas pelas n observações de cada uma das variáveis preditoras.

O vector \mathbf{b} dos $p+1$ parâmetros ajustados é dado pelo produto matricial do enunciado: $\mathbf{b} = (\mathbf{X}^t \mathbf{X})^{-1} (\mathbf{X}^t \mathbf{y})$. Um produto matricial no R é indicado pelo operador “%*%”, enquanto que uma inversa matricial é calculada pelo comando `solve`. A transposta duma matriz é dada pelo comando `t`. Logo, o vector \mathbf{b} obtém-se da seguinte forma:

```

> solve(t(X) %*% X) %*% t(X) %*% brix$Brix
      [,1]
(Intercept) 6.08877506
Diametro    1.27092840
Altura      -0.70967465
Peso        -0.20452522
pH          0.51556821
Acucar      0.08971091

```

Como se pode confirmar, trata-se dos valores já obtidos através do comando `lm`.

10. Começemos por recordar alguns resultados já previamente discutidos:

- Sabemos que, para qualquer conjunto de n pares de observações, se tem:

$$(n-1) \text{cov}_{xy} = \sum_{i=1}^n (x_i - \bar{x})y_i = \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} \Leftrightarrow \sum_{i=1}^n x_i y_i = (n-1) \text{cov}_{xy} + n \bar{x} \bar{y}. \quad (1)$$

- Tomando $y_i = x_i$, para todo o i , na fórmula anterior, obtém-se:

$$(n-1) s_x^2 = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n \bar{x}^2 \Leftrightarrow \sum_{i=1}^n x_i^2 = (n-1) s_x^2 + n \bar{x}^2 \quad (2)$$

- O produto de matrizes AB só é possível quando o número de colunas da matriz A for igual ao número de linhas da matriz B (matrizes *compatíveis* para a multiplicação). Se A é de dimensão $p \times q$ e B de dimensão $q \times r$, o produto AB é de dimensão $p \times r$.

- O elemento na linha i , coluna j , dum produto matricial AB , é dado pelo *produto interno* da linha i de A com a coluna j de B : $(AB)_{ij} = (a_{i1} \ a_{i2} \ \dots \ a_{iq}) \begin{pmatrix} b_{1j} \\ b_{2j} \\ \vdots \\ b_{qj} \end{pmatrix} = \sum_{k=1}^q a_{ik} b_{kj}$.

- O produto interno de dois vectores n -dimensionais \vec{x} e \vec{y} é dado por $\vec{x}^t \vec{y} = \sum_{i=1}^n x_i y_i$. No caso de um dos vectores ser o vector de n uns, $\mathbf{1}_n$, o produto interno resulta na soma dos elementos do outro vector, ou seja, em n vezes a média dos elementos do outro vector:

$$\mathbf{1}_n^t \vec{x} = \sum_{i=1}^n x_i = n \bar{x}.$$

- A *matriz inversa* dum matriz $n \times n$ A é definida (caso exista) como a matriz (única) A^{-1} , também de dimensão $n \times n$, tal que $AA^{-1} = \mathbf{I}_n$, onde \mathbf{I}_n é a matriz identidade de dimensão $n \times n$ (recorde-se que uma matriz identidade é uma matriz quadrada com todos os elementos diagonais iguais a 1 e todos os elementos não diagonais iguais a zero).

- No caso de A ser uma matriz 2×2 , de elementos $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$, a matriz inversa é dada (verifique!) por:

$$A^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \quad (3)$$

esta matriz inversa existe *se e só se* o *determinante* $ad - bc \neq 0$.

Com estes resultados prévios, as contas do exercício resultam de forma simples:

- (a) A matriz do modelo \mathbf{X} é de dimensão $n \times (p+1)$, que no caso dum regressão linear simples ($p=1$), significa $n \times 2$. Tem uma primeira coluna de uns (o vector $\mathbf{1}_n$) e uma segunda coluna com os n valores observados da variável preditora x , coluna essa que designamos pelo vector \vec{x} . Logo, a sua transposta \mathbf{X}^t é de dimensão $2 \times n$. Como o vector \vec{y} é de dimensão $n \times 1$, o produto $\mathbf{X}\vec{y}$ é possível e o resultado é um vector de dimensão 2×1 . O primeiro elemento (na posição (1,1)) desse produto é dada pelo produto interno da primeira linha de \mathbf{X}^t com a primeira e única coluna de \vec{y} , ou seja, por $\mathbf{1}_n^t \vec{y} = \sum_{i=1}^n y_i = n \bar{y}$. O segundo elemento (posição (2,1)) desse vector é dado pelo produto interno da segunda linha de \mathbf{X}^t e a única coluna de \vec{y} , ou seja, por $\vec{x}^t \vec{y} = \sum_{i=1}^n x_i y_i = (n-1) \text{cov}_{xy} + n \bar{x} \bar{y}$, tendo em conta a equação (1).
- (b) Tendo em conta que \mathbf{X}^t é de dimensão $2 \times n$ e \mathbf{X} é de dimensão $n \times 2$, o produto $\mathbf{X}^t \mathbf{X}$ é possível e de dimensão 2×2 . O elemento na posição (1,1) é o produto interno da primeira linha de \mathbf{X}^t ($\mathbf{1}_n$) com a primeira coluna de \mathbf{X} (igualmente $\mathbf{1}_n$), logo é: $\mathbf{1}_n^t \mathbf{1}_n = n$. O elemento na posição (1,2) é o produto interno da primeira linha de \mathbf{X}^t ($\mathbf{1}_n$) e segunda coluna de \mathbf{X} (\vec{x}), logo é $\mathbf{1}_n^t \vec{x} = \sum_{i=1}^n x_i = n \bar{x}$. O elemento na posição (2,1) é o produto interno da segunda linha de \mathbf{X}^t (\vec{x}) com a primeira coluna de \mathbf{X} ($\mathbf{1}_n$), logo é também $n \bar{x}$. Finalmente, o elemento na posição (2,2) é o produto interno da segunda linha de \mathbf{X}^t (\vec{x}) com a segunda coluna de \mathbf{X} (\vec{x}), ou seja, $\vec{x}^t \vec{x} = \sum_{i=1}^n x_i^2$. Fica assim provado o resultado do enunciado.

- (c) A primeira expressão da inversa dada no enunciado vem directamente de aplicar a fórmula (3) à matriz $(\mathbf{X}^t \mathbf{X})$ obtida na alínea anterior. Apenas há que confirmar a expressão do determinante $ad-bc = n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2 = n \sum_{i=1}^n x_i^2 - (n\bar{x})^2 = n \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) = n(n-1) s_x^2$, tendo em conta a fórmula (2). Igualmente a partir da fórmula (2) obtém-se a expressão alternativa do elemento na posição (1,1), que surge na segunda expressão para $(\mathbf{X}^t \mathbf{X})^{-1}$. Admitindo um contexto inferencial, ao multiplicar a matriz $(\mathbf{X}^t \mathbf{X})^{-1}$ pela variância σ^2 dos erros aleatórios obtém-se a matriz

$$\sigma^2(\mathbf{X}^t \mathbf{X})^{-1} = \begin{bmatrix} \sigma^2 \frac{(n-1)s_x^2 + n\bar{x}^2}{n(n-1)s_x^2} & \frac{-n\bar{x}\sigma^2}{n(n-1)s_x^2} \\ \frac{-n\bar{x}\sigma^2}{n(n-1)s_x^2} & \frac{n\sigma^2}{n(n-1)s_x^2} \end{bmatrix} = \begin{bmatrix} \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2} \right] & \frac{-\bar{x}\sigma^2}{(n-1)s_x^2} \\ \frac{-\bar{x}\sigma^2}{(n-1)s_x^2} & \frac{\sigma^2}{(n-1)s_x^2} \end{bmatrix}$$

No canto superior esquerdo tem-se a expressão de $V[\hat{\beta}_0]$. No canto inferior direito a expressão de $V[\hat{\beta}_1]$. O elemento comum às duas posições não diagonais é $Cov[\hat{\beta}_0, \hat{\beta}_1] = Cov[\hat{\beta}_1, \hat{\beta}_0]$.

- (d) Usando as expressões finais obtidas nas alíneas (c) e (a), obtém-se

$$\begin{aligned} (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \bar{\mathbf{y}} &= \frac{1}{n(n-1)s_x^2} \begin{bmatrix} (n-1)s_x^2 + n\bar{x}^2 & -n\bar{x} \\ -n\bar{x} & n \end{bmatrix} \begin{bmatrix} n\bar{y} \\ (n-1)cov_{xy} + n\bar{x}\bar{y} \end{bmatrix} \\ &= \frac{1}{n(n-1)s_x^2} \begin{bmatrix} (n-1)s_x^2 n\bar{y} + n^2\bar{x}^2\bar{y} - n\bar{x}(n-1)cov_{xy} - n^2\bar{x}^2\bar{y} \\ -n^2\bar{x}\bar{y} + n(n-1)cov_{xy} + n^2\bar{x}\bar{y} \end{bmatrix} \\ &= \begin{bmatrix} \frac{n(n-1)s_x^2\bar{y} - n(n-1)cov_{xy}\bar{x}}{n(n-1)s_x^2} \\ \frac{n(n-1)cov_{xy}}{n(n-1)s_x^2} \end{bmatrix} = \begin{bmatrix} \bar{y} - b_1\bar{x} \\ b_1 \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}. \end{aligned}$$

11. Sabemos que a matriz de projecção ortogonal referida é dada por $\mathbf{H} = \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t$, onde \mathbf{X} é a matriz do modelo, ou seja, a matriz de dimensões $n \times (p+1)$ que tem na primeira coluna, n uns, e em cada uma das p restantes colunas, as n observações de cada variável preditora. Ora,

- (a) A idempotência é fácil de verificar, tendo em conta que $(\mathbf{X}^t \mathbf{X})^{-1}$ é a matriz inversa de $\mathbf{X}^t \mathbf{X}$:

$$\mathbf{H} \mathbf{H} = \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t = \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t = \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t = \mathbf{H}.$$

A simetria resulta de três propriedades conhecidas de matrizes: a transposta dum matriz transposta é a matriz original $((\mathbf{A}^t)^t = \mathbf{A})$; a transposta dum produto de matrizes é o produto das correspondentes transpostas, pela ordem inversa $((\mathbf{A}\mathbf{B})^t = \mathbf{B}^t \mathbf{A}^t)$; e a transposta dum matriz inversa é a inversa da transposta $((\mathbf{A}^{-1})^t = (\mathbf{A}^t)^{-1})$. De facto, tem-se:

$$\mathbf{H}^t = [\mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t]^t = \mathbf{X}[(\mathbf{X}^t \mathbf{X})^{-1}]^t \mathbf{X}^t = \mathbf{X}[(\mathbf{X}^t \mathbf{X})^t]^{-1} \mathbf{X}^t = \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t = \mathbf{H}.$$

- (b) Como foi visto nas aulas, qualquer vector do subespaço das colunas da matriz \mathbf{X} , ou seja, do subespaço $\mathcal{C}(\mathbf{X}) \subset \mathbb{R}^n$, se pode escrever como $\mathbf{X}\mathbf{a}$, onde $\mathbf{a} \in \mathbb{R}^{p+1}$ é o vector dos $p+1$ coeficientes na combinação linear das colunas de \mathbf{X} . Ora, a projecção ortogonal deste vector sobre o subespaço $\mathcal{C}(\mathbf{X})$ (que já o contém) é dada por

$$\mathbf{H}\mathbf{X}\mathbf{a} = \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t (\mathbf{X}\mathbf{a}) = \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} (\mathbf{X}^t \mathbf{X}) \mathbf{a} = \mathbf{X}\mathbf{a}.$$

Assim, o vector $\mathbf{X}\mathbf{a}$ fica igual após a projecção.

- (c) Por definição, o vector dos valores ajustados é dado por $\vec{\hat{y}} = \mathbf{H}\vec{y}$. Ora, a média desses valores ajustados, que podemos representar por $\bar{\hat{y}} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i$, pode ser calculado tomando o produto interno do vector $\mathbf{1}_n$ de n uns com o vector $\vec{\hat{y}}$, uma vez que esse produto interno devolve a soma dos elementos de $\vec{\hat{y}}$. Assim, a média dos valores ajustados é $\bar{\hat{y}} = \frac{1}{n} \mathbf{1}_n^t \vec{\hat{y}} = \frac{1}{n} \mathbf{1}_n^t \mathbf{H}\vec{y} = \frac{1}{n} (\mathbf{H}\mathbf{1}_n)^t \vec{y} = \frac{1}{n} \mathbf{1}_n^t \vec{y}$, uma vez que $\mathbf{H}\mathbf{1}_n = \mathbf{1}_n$, já que a projecção ortogonal dum vector num subespaço onde ele já está contido deixa esse vector invariante, e o vector $\mathbf{1}_n$ pertence ao subespaço $\mathcal{C}(\mathbf{X})$ sobre o qual \mathbf{H} projecta, já que é a primeira das colunas da matriz \mathbf{X} . Mas a expressão final obtida, $\frac{1}{n} \mathbf{1}_n^t \vec{y}$ é a média \bar{y} dos valores observados de Y (já que $\mathbf{1}_n^t \vec{y}$ devolve a soma dos elementos do vector dessas observações, \vec{y}). Assim, na regressão linear múltipla, valores observados de Y e correspondentes valores ajustados partilham o mesmo valor médio.
- (d) O vector dos resíduos é dado por $\vec{e} = \vec{y} - \vec{\hat{y}} = \vec{y} - \mathbf{H}\vec{y}$. A soma dos resíduos resulta do produto interno do vector \vec{e} e o vector $\mathbf{1}_n$. Assim, tem-se (tendo também em conta a discussão das alíneas anteriores) $\mathbf{1}_n^t \vec{e} = \mathbf{1}_n^t (\vec{y} - \mathbf{H}\vec{y}) = \mathbf{1}_n^t \vec{y} - \mathbf{1}_n^t \mathbf{H}\vec{y} = \mathbf{1}_n^t \vec{y} - \mathbf{1}_n^t \vec{y} = 0$.

12. (a) A matriz de projecção ortogonal $\mathbf{P} = \mathbf{1}_n (\mathbf{1}_n^t \mathbf{1}_n)^{-1} \mathbf{1}_n^t$ é de dimensão $n \times n$ (confirme!), uma vez que o vector $\mathbf{1}_n$ é $n \times 1$. Mas o seu cálculo é facilitado pelo facto de que $\mathbf{1}_n^t \mathbf{1}_n$ é, neste caso, um escalar. Concretamente, $\mathbf{1}_n^t \mathbf{1}_n = n$, pelo que $(\mathbf{1}_n^t \mathbf{1}_n)^{-1} = \frac{1}{n}$. Logo $\mathbf{P} = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^t$. O produto $\mathbf{1}_n \mathbf{1}_n^t$ resulta numa matriz $n \times n$ com todos os elementos iguais a 1 (não confundir com o produto pela ordem inversa, $\mathbf{1}_n^t \mathbf{1}_n$: recorde-se que o produto de matrizes **não** é comutativo). Assim,

$$\mathbf{P} = \frac{1}{n} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}$$

- (b) A projecção ortogonal do vector $\vec{x} = (x_1, x_2, \dots, x_n)^t$ (cujos elementos serão por nós encarados como n observações duma variável X) sobre o subespaço gerado pelo vector dos uns $\mathbf{1}_n$ é:

$$\mathbf{P}\vec{x} = \frac{1}{n} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \bar{x} \\ \bar{x} \\ \bar{x} \\ \vdots \\ \bar{x} \end{bmatrix} = \bar{x} \cdot \mathbf{1}_n$$

onde $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ é a média dos valores do vector \vec{x} . Ou seja, o vector projectado é um múltiplo escalar do vector dos uns (como são todos os vectores que pertencem a $\mathcal{C}(\mathbf{1}_n)$, uma vez que as combinações lineares dum qualquer vector são sempre múltiplos escalares desse vector). Mas a constante de multiplicação desse vector projectado tem significado estatístico: é a média dos valores do vector \vec{x} .

- (c) É característico da matriz identidade \mathbf{I} que, para qualquer vector \vec{x} se tem $\mathbf{I}\vec{x} = \vec{x}$. Logo,

tendo em conta o resultado da alínea anterior, tem-se:

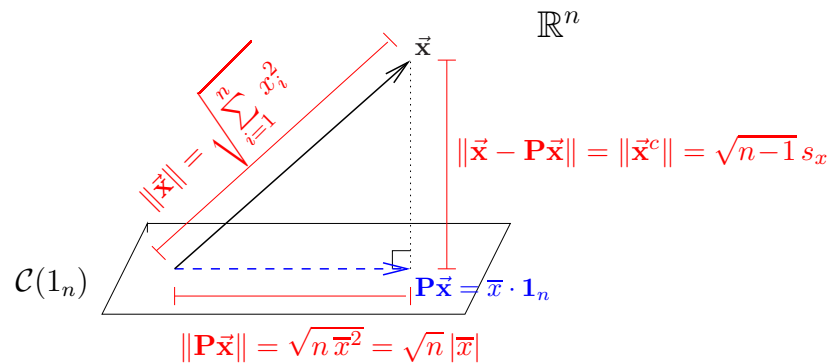
$$(\mathbf{I} - \mathbf{P})\vec{x} = \mathbf{I}\vec{x} - \mathbf{P}\vec{x} = \vec{x} - \mathbf{P}\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} - \begin{bmatrix} \bar{x} \\ \bar{x} \\ \bar{x} \\ \vdots \\ \bar{x} \end{bmatrix} = \begin{bmatrix} x_1 - \bar{x} \\ x_2 - \bar{x} \\ x_3 - \bar{x} \\ \vdots \\ x_n - \bar{x} \end{bmatrix} = \vec{x}^c$$

- (d) A norma do vector \vec{x}^c é, por definição, a raiz quadrada da soma dos quadrados dos seus elementos. Logo, tendo em conta a natureza dos elementos do vector \vec{x}^c (ver a alínea anterior), tem-se:

$$\|\vec{x}^c\| = \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{(n-1)s_x^2} = \sqrt{n-1} s_x,$$

ou seja, a norma é proporcional ao desvio padrão s_x dos valores do vector \vec{x} (sendo a constante de proporcionalidade $\sqrt{n-1}$).

- (e) A situação considerada nas alíneas anteriores tem a seguinte representação gráfica:



Nota: O subespaço $\mathcal{C}(1_n)$ é gerado por um único vector, $\mathbf{1}_n$, pelo que em termos geométricos é uma linha recta que atravessa a origem (um subespaço de dimensão 1). Esse subespaço foi representado aqui por um plano para manter coerência com as representações gráficas usadas nas aulas, salientando que se trata do mesmo conceito de projecções ortogonais.

Aplicando o Teorema de Pitágoras ao triângulo rectângulo indicado, tem-se:

$$\sum_{i=1}^n x_i^2 = (n-1)s_x^2 + n\bar{x}^2 \Leftrightarrow s_x^2 = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right),$$

que é a chamada *fórmula computacional da variância*.

13. Note-se que a matriz \mathbf{P}_{1_n} referida neste exercício (e que será representada apenas por \mathbf{P} no que se segue) é a mesma que foi discutida no Exercício 12. Assim, o vector $\vec{y} - \mathbf{P}\vec{y}$ é o vector centrado das observações de \vec{y} :

$$\vec{y} - \mathbf{P}\vec{y} = \begin{bmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \\ y_3 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{bmatrix} = \vec{y}^c$$

A norma deste vector, ao quadrado, é a soma dos quadrados dos seus elementos, ou seja, $SQT = \sum_{i=1}^n (y_i - \bar{y})^2$. De forma análoga, e como o vector $\hat{\mathbf{y}}$ dos valores ajustados é dado por $\hat{\mathbf{y}} = \mathbf{X}\vec{\beta} = \mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t\vec{y} = \mathbf{H}\vec{y}$, temos que o vector $\mathbf{H}\vec{y} - \mathbf{P}\vec{y}$ tem como elementos $\hat{y}_i - \bar{y}$:

$$\mathbf{H}\vec{y} - \mathbf{P}\vec{y} = \begin{bmatrix} \hat{y}_1 - \bar{y} \\ \hat{y}_2 - \bar{y} \\ \hat{y}_3 - \bar{y} \\ \vdots \\ \hat{y}_n - \bar{y} \end{bmatrix}$$

pelo que o quadrado da sua norma é $SQR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$. Finalmente, o vector $\vec{y} - \mathbf{H}\vec{y} = \vec{y} - \hat{\mathbf{y}}$ é o vector dos resíduos, e a sua norma ao quadrado é $SQRE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$.

Nas aulas viu-se geometricamente que o Teorema de Pitágoras garante que $SQT = SQR + SQRE$. Neste exercício pede-se para confirmar tal facto do ponto de vista algébrico. Tendo em conta que as Somas de Quadrados são os quadrados das normas acima indicados, e recordando as propriedades de normas, temos:

$$\begin{aligned} SQT &= \|\vec{y} - \mathbf{P}\vec{y}\|^2 = \|(\vec{y} - \mathbf{H}\vec{y}) + (\mathbf{H}\vec{y} - \mathbf{P}\vec{y})\|^2 \\ &= \|\vec{y} - \mathbf{H}\vec{y}\|^2 + \|\mathbf{H}\vec{y} - \mathbf{P}\vec{y}\|^2 + 2(\vec{y} - \mathbf{H}\vec{y})|(\mathbf{H}\vec{y} - \mathbf{P}\vec{y}) \\ &= SQR + SQRE + 2(\vec{y} - \mathbf{H}\vec{y})|(\mathbf{H}\vec{y} - \mathbf{P}\vec{y}) \end{aligned}$$

onde na última parcela surge o produto interno entre os vectores $\vec{y} - \mathbf{H}\vec{y}$ e $\mathbf{H}\vec{y} - \mathbf{P}\vec{y}$. Este produto interno tem de ser nulo, para ser verdade a relação entre as Somas de Quadrados. Ora,

$$\begin{aligned} (\vec{y} - \mathbf{H}\vec{y})|(\mathbf{H}\vec{y} - \mathbf{P}\vec{y}) &= (\vec{y} - \mathbf{H}\vec{y})^t(\mathbf{H}\vec{y} - \mathbf{P}\vec{y}) \\ &= \vec{y}^t\mathbf{H}\vec{y} - \vec{y}^t\mathbf{P}\vec{y} - (\mathbf{H}\vec{y})^t\mathbf{H}\vec{y} + (\mathbf{H}\vec{y})^t\mathbf{P}\vec{y} \\ &= \vec{y}^t\mathbf{H}\vec{y} - \vec{y}^t\mathbf{P}\vec{y} - \vec{y}^t\mathbf{H}^t\mathbf{H}\vec{y} + \vec{y}^t\mathbf{H}^t\mathbf{P}\vec{y}, \end{aligned} \quad (4)$$

tendo em conta que, em qualquer produto matricial, a transposta do produto é o produto das transpostas pela ordem inversa ($(\mathbf{AB})^t = \mathbf{B}^t\mathbf{A}^t$). Mas (tal como se viu no Exercício 11) \mathbf{H} é uma matriz simétrica: $\mathbf{H}^t = [\mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t]^t = \mathbf{X}[(\mathbf{X}^t\mathbf{X})^{-1}]^t\mathbf{X}^t = \mathbf{X}[(\mathbf{X}^t\mathbf{X})^t]^{-1}\mathbf{X}^t = \mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t = \mathbf{H}$, tendo em conta que, para qualquer matriz invertível, a inversa da transposta é a transposta da inversa ($(\mathbf{A}^t)^{-1} = (\mathbf{A}^{-1})^t$), e que a transposta duma transposta é a matriz original ($(\mathbf{A}^t)^t = \mathbf{A}$). Por outro lado, $\mathbf{H}\mathbf{H} = \mathbf{H}$, porque $\mathbf{H}\mathbf{H} = [\mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t][\mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t] = \mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}(\mathbf{X}^t\mathbf{X})(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t = \mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t = \mathbf{H}$. Logo, a terceira parcela na equação (4) vem igual à primeira ($\vec{y}^t\mathbf{H}\vec{y}$), mas de sinal contrário, cancelando. Por seu lado, e de novo usando a simetria de \mathbf{H} , a matriz da última parcela em (4) vem $\mathbf{H}^t\mathbf{P} = \mathbf{H}\mathbf{P} = \mathbf{H}\mathbf{1}_n(\mathbf{1}_n^t\mathbf{1}_n)^{-1}\mathbf{1}_n^t$. Mas (como se viu nas aulas teóricas) $\mathbf{H}\mathbf{1}_n = \mathbf{1}_n$, uma vez que o vector $\mathbf{1}_n$ pertence ao subespaço $\mathcal{C}(\mathbf{X})$ sobre o qual a matriz \mathbf{H} projecta, e qualquer vector fica invariante quando projectado sobre um subespaço ao qual pertence. Logo, $\mathbf{H}\mathbf{P} = \mathbf{1}_n(\mathbf{1}_n^t\mathbf{1}_n)^{-1}\mathbf{1}_n^t = \mathbf{P}$. Assim, a última parcela da equação (4) vem igual à segunda ($\vec{y}^t\mathbf{P}\vec{y}$), mas com sinal trocado, pelo que essas duas parcelas também cancelam e o produto interno indicado nessa equação anula-se.

14. A informação essencial sobre a regressão pedida pode ser obtida através do comando `summary`:

```

> iris.lm <- lm(Petal.Width ~ Petal.Length, data=iris)
> summary(iris.lm)
Call: lm(formula = Petal.Width ~ Petal.Length, data = iris)
(...)
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.363076  0.039762  -9.131  4.7e-16 ***
Petal.Length  0.415755  0.009582  43.387  < 2e-16 ***
(...)
Residual standard error: 0.2065 on 148 degrees of freedom
Multiple R-squared:  0.9271, Adjusted R-squared:  0.9266
F-statistic: 1882 on 1 and 148 DF,  p-value: < 2.2e-16

```

- (a) As estimativas dos desvios padrão associados à estimação de cada parâmetro são indicadas na coluna de nome **Std.Error** (erro padrão). O desvio padrão associado à estimação da ordenada na origem é $\hat{\sigma}_{\hat{\beta}_0} = 0.039762$. A variância correspondente é o quadrado deste valor, $\hat{\sigma}_{\hat{\beta}_0}^2 = 0.001581$. É também possível calcular esta variância estimada a partir da sua fórmula: $\hat{\sigma}_{\hat{\beta}_0}^2 = QMRE (\mathbf{X}^t \mathbf{X})_{(1,1)}^{-1}$. Tratando-se duma regressão linear *simples*, é possível provar a seguinte fórmula alternativa: $\hat{\sigma}_{\hat{\beta}_0}^2 = QMRE \left[\frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2} \right]$, onde \bar{x} e s_x^2 indicam, respectivamente, a média e a variância amostral dos n valores de X observados. O valor de $QMRE$ pode ser obtido a partir da listagem acima: **Residual standard error** indica o valor $\sqrt{QMRE} = 0.2065$. Os outros valores constantes da expressão podem ser calculados como em exercícios anteriores. O desvio padrão associado à estimação do declive da recta é $\hat{\sigma}_{\hat{\beta}_1} = 0.009582$, e o seu quadrado é a variância estimada de $\hat{\beta}_1$: $\hat{\sigma}_{\hat{\beta}_1}^2 = 9.181472 \times 10^{-5}$. Este valor pode ser obtido a partir da expressão $\hat{\sigma}_{\hat{\beta}_1}^2 = QMRE \cdot (\mathbf{X}^t \mathbf{X})_{(2,2)}^{-1}$. Também neste caso, e tratando-se duma regressão linear simples, tem-se a expressão alternativa: $\hat{\sigma}_{\hat{\beta}_1}^2 = \frac{QMRE}{(n-1)s_x^2}$.
- (b) Um intervalo a $(1 - \alpha) \times 100\%$ de confiança para β_1 é: $\left] b_1 - t_{\frac{\alpha}{2}(n-2)} \hat{\sigma}_{\hat{\beta}_1}, b_1 + t_{\frac{\alpha}{2}(n-2)} \hat{\sigma}_{\hat{\beta}_1} \right[$, sendo neste caso $\alpha = 0.05$, $n = 150$, $b_1 = 0.415755$, $\hat{\sigma}_{\hat{\beta}_1} = 0.009582$ e $t_{0.025(148)} = 1.976122$. Logo, o IC a 95% de confiança para o declive da recta é $] 0.39682, 0.43469 [$. Esta é a gama de valores admissíveis (a 95% de confiança) para o declive da recta relacionando largura e comprimento das pétalas dos lírios (das três espécies analisadas). Os intervalos de confiança dos dois parâmetros da recta podem ser obtidos no R através do comando:

```

> confint(iris.lm)
                2.5 %      97.5 %
(Intercept) -0.4416501 -0.2845010
Petal.Length  0.3968193  0.4346915

```

- (c) Analogamente, um IC a $(1 - \alpha) \times 100\%$ de confiança para β_0 é:

$$\left] b_0 - t_{\frac{\alpha}{2}(n-2)} \cdot \hat{\sigma}_{\hat{\beta}_0}, b_0 + t_{\frac{\alpha}{2}(n-2)} \cdot \hat{\sigma}_{\hat{\beta}_0} \right[$$

Neste exemplo, $b_0 = -0.363076$ e $\hat{\sigma}_{\hat{\beta}_0} = 0.039762$. O valor tabelado da distribuição t , para um intervalo a 95% de confiança, é o mesmo que na alínea anterior: $t_{0.025(148)} = 1.976122$. Logo, o intervalo de confiança pedido é $] -0.4416501, -0.2845010 [$. Repare-se na maior amplitude deste intervalo, em relação ao IC para o declive populacional β_1 , o que é consequência directa da maior variabilidade associada à estimação de β_0 (o valor de $\hat{\sigma}_{\hat{\beta}_0}$ é cerca de 4 vezes o valor de $\hat{\sigma}_{\hat{\beta}_1}$). A partir das fórmulas para estes dois erros padrão, é possível verificar que este

maior valor de $\hat{\sigma}_{\hat{\beta}_0}$ resulta, não tanto da parcela adicional $\frac{1}{n}$ (como $n = 150$, esta parcela é pequena) mas sobretudo do \bar{x}^2 que surge no numerador da segunda parcela. De facto, a média das observações do comprimento de pétalas é aproximadamente $\bar{x} = 3.758$.

- (d) A frase do enunciado traduz-se por “ $\beta_1 = 0.5$ ”. Assim, faremos um teste de hipóteses desta hipótese nula, contra a hipótese alternativa $H_1 : \beta_1 \neq 0.5$. Os cinco passos do teste são:

Hipóteses: $H_0 : \beta_1 = 0.5$ vs. $H_1 : \beta_1 \neq 0.5$.

Estatística do teste: $T = \frac{\hat{\beta}_1 - \beta_{1|H_0}}{\hat{\sigma}_{\hat{\beta}_1}} \cap t_{n-2}$

Nível de significância: $\alpha = 0.05$.

Região Crítica (Bilateral): Rejeitar H_0 se $|T_{calc}| > t_{\frac{\alpha}{2}(n-2)} = t_{0.025(148)} = 1.976122$.

Conclusões: O valor calculado da estatística do teste é: $T_{calc} = \frac{0.415755 - 0.5}{0.009582} = -8.792006$.

Logo, rejeita-se claramente a hipótese nula que por cada centímetro a mais no comprimento da pétala, é de esperar meio centímetro a mais na largura da pétala.

- (e) A hipótese referida no enunciado é que $\beta_1 < 0.5$. Neste caso, a opção entre colocar esta hipótese em H_0 ou em H_1 corresponde à opção entre dar, ou não, o benefício da dúvida a esta hipótese. Seja como for, o valor de fronteira (0.5) terá de pertencer à hipótese nula. Vamos optar por *não* dar o benefício da dúvida à hipótese indicada no enunciado:

Hipóteses: $H_0 : \beta_1 \geq 0.5$ vs. $H_1 : \beta_1 < 0.5$.

Estatística do teste: $T = \frac{\hat{\beta}_1 - 0.5}{\hat{\sigma}_{\hat{\beta}_1}} \cap t_{n-2}$

Nível de significância: $\alpha = 0.05$.

Região Crítica (Unilateral esquerda): Rej. H_0 se $T_{calc} < -t_{\alpha(n-2)} = -t_{0.05(148)} = -1.655215$.

Conclusões: O valor calculado da estatística do teste é igual ao da alínea anterior: $T_{calc} = \frac{0.415755 - 0.5}{0.009582} = -8.792006$. Logo, rejeita-se a hipótese nula, optando-se por H_1 . Pode afirmar-se que é estatisticamente significativa a conclusão que, por cada centímetro a mais no comprimento da pétala, em média a respectiva largura cresce menos do que 0.5cm.

- (f) A afirmação do enunciado corresponde à hipótese $\beta_1 = 0$. De facto, se $\beta_1 = 0$, a equação do modelo que relaciona x e Y reduz-se a $Y_i = \beta_0 + \epsilon_i$, não existindo relação linear entre x e Y . O teste às hipóteses $H_0 : \beta_1 = 0$ vs. $H_1 : \beta_1 \neq 0$ pode ser feito como na alínea 14d) acima. No entanto, para o caso particular do valor do parâmetro $\beta_1 = 0$ a informação relativa ao teste já é indicada na listagem produzida pelo comando `summary`, nas terceira e quarta colunas da tabela `Coefficients`. Neste caso, o valor calculado da estatística é $T_{calc} = \frac{0.4157550}{0.009582} = 43.387$. Tendo em conta que a região crítica é igual à da alínea 14d), tem-se uma rejeição clara da hipótese nula $\beta_1 = 0$: o valor estimado $b_1 = 0.415755$ é *significativamente diferente* de zero (ao nível $\alpha = 0.05$), pelo que a recta tem alguma utilidade para prever valores de y (largura da pétala) a partir dos valores de x (comprimento da pétala). Esta conclusão também se pode justificar a partir do valor de prova (p -value) do valor calculado da estatística, que é muito pequeno, sendo mesmo inferior à precisão de máquina, $p < 2 \times 10^{-16}$. Mesmo para níveis de significância como $\alpha = 0.01$ ou $\alpha = 0.005$, a conclusão seria a de rejeição de H_0 .
- (g) Uma abordagem alternativa para a questão estudada na alínea anterior será a de efectuar um teste de ajustamento global (teste F) à regressão ajustada. No nosso caso, e definindo \mathcal{R}^2 como o coeficiente de determinação populacional, tem-se:

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

Estatística do teste: $F = \frac{QMR}{QMRE} = (n-2) \frac{R^2}{1-R^2} \cap F_{(1,n-2)}$, sob H_0 .

Nível de significância: $\alpha = 0.05$.

Região Crítica (Unilateral direita): Rej. H_0 se $F_{calc} > f_{\alpha(1,n-2)} = f_{0.05(1,148)} = 3.905$.

Conclusões: O valor calculado da estatística é: $F_{calc} = 148 \times \frac{0.9271}{1-0.9271} = 1882.178$. Logo, rejeita-se claramente a hipótese nula, que corresponde à hipótese dum ajustamento inútil do modelo. A resposta é coerente com a alínea anterior.

NOTA: Repare-se que o comando `summary` do R, quando aplicado ao ajustamento dum regressão, indica na última linha das listagens o valor da estatística calculada F_{calc} , os respectivos graus de liberdade associados, e o valor de prova (*p-value*) correspondente.

- (h) A largura esperada dum pétala cujo comprimento seja $x = 4.5\text{cm}$ é dada por $\hat{\mu} = b_0 + b_1 4.5 = -0.363076 + 0.415755 \times 4.5 = 1.507821$. No R, este resultado pode ser obtido através do comando `predict`:

```
> predict(iris.lm, new=data.frame(Petal.Length=4.5))
1
1.507824
```

O intervalo de confiança para $\mu_{x=4.5} = E[Y|X = 4.5]$ é dado por:

$$\left[(b_0 + b_1 x) - t_{\frac{\alpha}{2}; n-2} \sqrt{QMRE \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{(n-1) s_x^2} \right]}, (b_0 + b_1 x) + t_{\frac{\alpha}{2}; n-2} \sqrt{QMRE \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{(n-1) s_x^2} \right]} \right]$$

em que $\hat{\mu} = b_0 + b_1 4.5 = 1.507821$, $t_{\frac{\alpha}{2}; n-2} = t_{0.025, 148} = 1.976122$, $QMRE = 0.2065^2$ (a partir da listagem acima dada). Por outro lado, a média e variância das $n = 150$ observações do preditor `Petal.Length` podem ser calculadas e resultam ser $\bar{x} = 3.758$ e $s_x^2 = 3.116278$. Assim, a 95% de confiança, o verdadeiro valor de $\mu_{x=4.5} = E[Y|X = 4.5]$ faz parte do intervalo $] 1.471666, 1.543982 [$. No R este intervalo de confiança pode ser obtido através do comando

```
> predict(iris.lm, new=data.frame(Petal.Length=4.5), int="conf")
      fit      lwr      upr
1 1.507824 1.471666 1.543982
```

Os extremos do intervalo são dados pelos valores `lwr` (de *lower*) e `upr` (de *upper*).

- (i) O intervalo *de predição* para o valor da variável resposta y (largura da pétala) associada a uma observação com $x = 4.5$ é dado por:

$$\left[(b_0 + b_1 x) - t_{\frac{\alpha}{2}; n-2} \sqrt{QMRE \left[1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{(n-1) s_x^2} \right]}, (b_0 + b_1 x) + t_{\frac{\alpha}{2}; n-2} \sqrt{QMRE \left[1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{(n-1) s_x^2} \right]} \right]$$

Em relação ao intervalo de confiança pedido na alínea anterior, apenas muda a expressão debaixo da raiz quadrada. No R este tipo de intervalo obtém-se com um comando muito semelhante ao anterior:

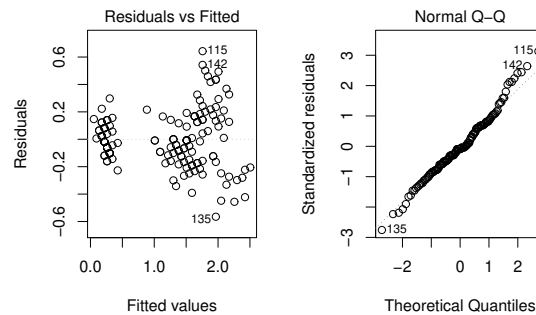
```
> predict(iris.lm, new=data.frame(Petal.Length=4.5), int="pred")
      fit      lwr      upr
1 1.507824 1.098187 1.917461
```

Como seria de esperar, trata-se dum intervalo bastante mais amplo: $] 1.098187, 1.917461 [$.

- (j) Dos gráficos de resíduos produzidos pelo comando

```
> plot(lm(Petal.Width ~ Petal.Length, data=iris), which=c(1,2))
```

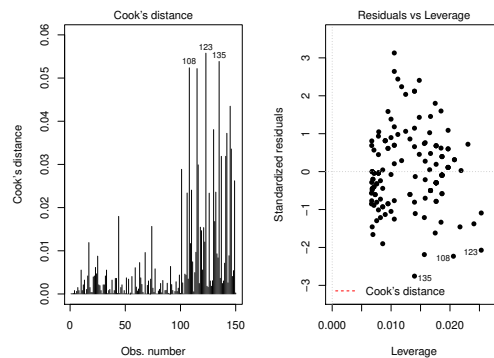
verifica-se que pode existir um problema em relação à hipótese de homogeneidade de variâncias. O gráfico da esquerda sugere que os lírios com comprimento de pétala mais pequeno (do lado esquerdo do gráfico) parecem ter menor variabilidade dos resíduos do que os restantes. Já a linearidade aproximada no *qq-plot* (gráfico da direita) não indicia a existência de problemas com a hipótese de normalidade.



Quanto aos gráficos de diagnóstico produzidos pelo comando

```
> plot(lm(Petal.Width ~ Petal.Length, data=iris),which=c(4,5))
```

observa-se no diagrama de barras das distâncias de Cook que, apesar de haver alguma variabilidade nos valores, em nenhum caso a distância de Cook excede o valor (bastante baixo) de 0.06. Assim, nenhuma observação se deve considerar influente. De igual forma, não há valores elevados do efeito alavanca (*leverage*), sendo o maior valor de h_{ii} inferior a 0.03 (ver o eixo horizontal do gráfico da direita). Assim, nenhuma observação se destaca por ter um efeito alavanca elevado.



- (k) Nas três subalíneas, as transformações de uma ou ambas as variáveis são transformações afins (lineares), razão pela qual o quadrado do coeficiente de correlação, ou seja, o coeficiente de determinação R^2 não sofre alteração. O que pode mudar são os parâmetros da recta de regressão ajustada.

- i. Neste caso, apenas a variável preditora sofre uma transformação multiplicativa, da forma $x \rightarrow x^* = cx$ (com $c = 10$). Vejamos qual o efeito deste tipo de transformações nos parâmetros da recta de regressão. Utilizando a habitual notação dos asteriscos para indicar os valores correspondentes à transformação, temos (tendo em conta que

$var(cx) = c^2 var(x)$:

$$b_1^* = \frac{cov_{x^*y}}{s_{x^*}^2} = \frac{cov(cx,y)}{c^2 s_x^2} = \frac{1}{c} \frac{cov(x,y)}{s_x^2} = \frac{1}{c} b_1 ;$$

e (tendo em conta o efeito de constantes multiplicativas sobre a média, ou seja, $\overline{x^*} = c\overline{x}$):

$$b_0^* = \overline{y} - b_1^* \overline{x^*} = \overline{y} - \frac{1}{c} b_1 \cdot c\overline{x} = \overline{y} - b_1 \overline{x} = b_0 .$$

Ou seja, neste caso a ordenada na origem não se altera, enquanto que o declive vem multiplicado por $\frac{1}{10}$. Confirmemos estes resultados com recurso ao R:

```
> lm(formula = Petal.Width ~ I(Petal.Length*10), data = iris)
Call:
lm(formula = Petal.Width ~ I(Petal.Length * 10), data = iris)
Coefficients:
```

```
  (Intercept)  I(Petal.Length * 10)
    -0.36308         0.04158
```

- ii. Neste caso, estamos perante uma transformação idêntica à usada na alínea 1i), pelo que já sabemos que iremos encontrar, quer a ordenada na origem, quer o declive, multiplicados por $c = 10$. Confirmando no R:

```
> lm(formula = I(Petal.Width*10) ~ Petal.Length, data = iris)
Call:
lm(formula = I(Petal.Width * 10) ~ Petal.Length, data = iris)
Coefficients:
```

```
  (Intercept)  Petal.Length
    -3.631         4.158
```

- iii. Finalmente, na conjugação das duas transformações discutidas nas subalíneas anteriores, e generalizando para as transformações multiplicativas $x \rightarrow cx$ e $y \rightarrow dy$, vem:

$$b_1^* = \frac{cov_{x^*y^*}}{s_{x^*}^2} = \frac{cov(cx,dy)}{c^2 s_x^2} = \frac{cd}{c^2} \frac{cov(x,y)}{s_x^2} = \frac{d}{c} b_1 ;$$

e:

$$b_0^* = \overline{y^*} - b_1^* \overline{x^*} = d\overline{y} - \frac{d}{c} b_1 \cdot c\overline{x} = d(\overline{y} - b_1 \overline{x}) = db_0 .$$

Como no nosso caso $c = d = 10$, o declive não se deve alterar, enquanto a ordenada na origem deverá ser 10 vezes maior do que no caso original dos dados não transformados.

```
> lm(formula = I(Petal.Width*10) ~ I(Petal.Length*10), data = iris)
```

```
Call:
```

```
lm(formula = I(Petal.Width * 10) ~ I(Petal.Length * 10), data = iris)
```

```
Coefficients:
```

```
  (Intercept)  I(Petal.Length * 10)
    -3.6308         0.4158
```

15. Seja $\mathbf{Z} = (Z_1, Z_2, \dots, Z_k)^t$. Tendo em conta a definição de vector esperado e de matriz de variâncias-covariâncias, bem como as propriedades dos valores esperados, variâncias e covariâncias de variáveis aleatórias (unidimensionais) tem-se:

(a)

$$E[\alpha \mathbf{Z}] = \begin{bmatrix} E[\alpha Z_1] \\ E[\alpha Z_2] \\ \vdots \\ E[\alpha Z_k] \end{bmatrix} = \begin{bmatrix} \alpha E[Z_1] \\ \alpha E[Z_2] \\ \vdots \\ \alpha E[Z_k] \end{bmatrix} = \alpha E[\mathbf{Z}] .$$

(b)

$$E[\mathbf{Z} + \mathbf{a}] = \begin{bmatrix} E[Z_1 + a_1] \\ E[Z_2 + a_2] \\ \vdots \\ E[Z_k + a_k] \end{bmatrix} = \begin{bmatrix} E[Z_1] + a_1 \\ E[Z_2] + a_2 \\ \vdots \\ E[Z_k] + a_k \end{bmatrix} = \begin{bmatrix} E[Z_1] \\ E[Z_2] \\ \vdots \\ E[Z_k] \end{bmatrix} + \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{bmatrix} = E[\mathbf{Z}] + \mathbf{a} .$$

(c)

$$\begin{aligned} V[\alpha \mathbf{Z}] &= \begin{bmatrix} V[\alpha Z_1] & Cov[\alpha Z_1, \alpha Z_2] & \cdots & Cov[\alpha Z_1, \alpha Z_k] \\ Cov[\alpha Z_2, \alpha Z_1] & V[\alpha Z_2] & \cdots & Cov[\alpha Z_2, \alpha Z_k] \\ \vdots & \vdots & \ddots & \vdots \\ Cov[\alpha Z_k, \alpha Z_1] & Cov[\alpha Z_k, \alpha Z_2] & \cdots & V[\alpha Z_k] \end{bmatrix} \\ &= \begin{bmatrix} \alpha^2 V[Z_1] & \alpha^2 Cov[Z_1, Z_2] & \cdots & \alpha^2 Cov[Z_1, Z_k] \\ \alpha^2 Cov[Z_2, Z_1] & \alpha^2 V[Z_2] & \cdots & \alpha^2 Cov[Z_2, Z_k] \\ \vdots & \vdots & \ddots & \vdots \\ \alpha^2 Cov[Z_k, Z_1] & \alpha^2 Cov[Z_k, Z_2] & \cdots & \alpha^2 V[Z_k] \end{bmatrix} = \alpha^2 V[\mathbf{Z}] \end{aligned}$$

(d)

$$\begin{aligned} V[\mathbf{Z} + \mathbf{a}] &= \begin{bmatrix} V[Z_1 + a_1] & Cov[Z_1 + a_1, Z_2 + a_2] & \cdots & Cov[Z_1 + a_1, Z_k + a_k] \\ Cov[Z_2 + a_2, Z_1 + a_1] & V[Z_2 + a_2] & \cdots & Cov[Z_2 + a_2, Z_k + a_k] \\ \vdots & \vdots & \ddots & \vdots \\ Cov[Z_k + a_k, Z_1 + a_1] & Cov[Z_k + a_k, Z_2 + a_2] & \cdots & V[Z_k + a_k] \end{bmatrix} \\ &= \begin{bmatrix} V[Z_1] & Cov[Z_1, Z_2] & \cdots & Cov[Z_1, Z_k] \\ Cov[Z_2, Z_1] & V[Z_2] & \cdots & Cov[Z_2, Z_k] \\ \vdots & \vdots & \ddots & \vdots \\ Cov[Z_k, Z_1] & Cov[Z_k, Z_2] & \cdots & V[Z_k] \end{bmatrix} = V[\mathbf{Z}] \end{aligned}$$

(e)

$$E[\mathbf{Z} + \vec{\mathbf{U}}] = \begin{bmatrix} E[Z_1 + U_1] \\ E[Z_2 + U_2] \\ \vdots \\ E[Z_k + U_k] \end{bmatrix} = \begin{bmatrix} E[Z_1] + E[U_1] \\ E[Z_2] + E[U_2] \\ \vdots \\ E[Z_k] + E[U_k] \end{bmatrix} = \begin{bmatrix} E[Z_1] \\ E[Z_2] \\ \vdots \\ E[Z_k] \end{bmatrix} + \begin{bmatrix} E[U_1] \\ E[U_2] \\ \vdots \\ E[U_k] \end{bmatrix} = E[\mathbf{Z}] + E[\vec{\mathbf{U}}] .$$

16. (a) Tem-se, recordando que $SQRE = SQT - SQR$,

$$F = \frac{QMR}{QMRE} = \frac{SQR/1}{SQRE/(n-2)} = (n-2) \frac{SQR}{SQT - SQR} = (n-2) \frac{R^2}{1 - R^2} ,$$

onde a última passagem resulta de dividir numerador e denominador por SQT .

(b) Como R^2 está entre 0 e 1, qualquer aumento de R^2 aumenta o numerador e diminui o denominador, provocando um aumento da fracção. Assim, a maiores valores de R^2 correspondem maiores valores da estatística F . Uma vez que o teste F tem hipótese nula $H_0 : \mathcal{R}^2 = 0$, é natural que se defina uma região crítica unilateral direita.

17. (a) Após descarregar o ficheiro `TN025.csv` para a sua área de trabalho, dê o comando:

```
> TN025 <- read.csv("TN025.csv")
```

Inspeccione a *data frame* `TN025` como indicado no enunciado.

- (b) O melhor preditor de P será a variável mais fortemente correlacionada com P. Apenas faz sentido calcular correlações entre variáveis *numéricas*, ou seja, as colunas 12 a 25. Eis as correlações (arredondadas a 2 casas decimais):

```
> round(cor(TN025[,12:25]), d=2)
```

	P	K	CA	MG	MN	CU	B	ZN	AL	FE	NA.	C	N	NKJ
P	1.00	0.80	0.48	0.88	0.63	0.47	0.62	0.47	0.32	0.16	0.28	-0.35	0.94	-0.42
K	0.80	1.00	0.33	0.83	0.37	0.62	0.63	0.43	0.32	0.14	0.53	-0.47	0.84	-0.36
CA	0.48	0.33	1.00	0.57	0.35	0.19	0.53	0.73	-0.07	-0.15	0.30	-0.39	0.41	0.31
MG	0.88	0.83	0.57	1.00	0.46	0.57	0.73	0.54	0.19	0.15	0.47	-0.51	0.90	-0.32
MN	0.63	0.37	0.35	0.46	1.00	0.02	0.49	0.22	0.36	0.06	0.03	-0.03	0.54	-0.32
CU	0.47	0.62	0.19	0.57	0.02	1.00	0.31	0.43	0.11	0.11	0.37	-0.37	0.48	-0.05
B	0.62	0.63	0.53	0.73	0.49	0.31	1.00	0.43	0.34	-0.01	0.42	-0.30	0.68	-0.18
ZN	0.47	0.43	0.73	0.54	0.22	0.43	0.43	1.00	-0.08	-0.04	0.52	-0.48	0.40	0.32
AL	0.32	0.32	-0.07	0.19	0.36	0.11	0.34	-0.08	1.00	-0.02	0.21	0.26	0.27	-0.16
FE	0.16	0.14	-0.15	0.15	0.06	0.11	-0.01	-0.04	-0.02	1.00	0.13	-0.19	0.20	-0.38
NA.	0.28	0.53	0.30	0.47	0.03	0.37	0.42	0.52	0.21	0.13	1.00	-0.42	0.30	0.03
C	-0.35	-0.47	-0.39	-0.51	-0.03	-0.37	-0.30	-0.48	0.26	-0.19	-0.42	1.00	-0.37	0.02
N	0.94	0.84	0.41	0.90	0.54	0.48	0.68	0.40	0.27	0.20	0.30	-0.37	1.00	-0.57
NKJ	-0.42	-0.36	0.31	-0.32	-0.32	-0.05	-0.18	0.32	-0.16	-0.38	0.03	0.02	-0.57	1.00

A variável mais fortemente correlacionada com o teor de fósforo é N, o teor de azoto.

- i. Eis os resultados relativos à regressão linear simples de P sobre N. Aparentemente estamos perante um bom modelo, que explica quase 88% da variabilidade dos teores de fósforo observados.

```
> TNP.lm <- lm(P ~ N , data=TN025)
> summary(TNP.lm)
[...]
```

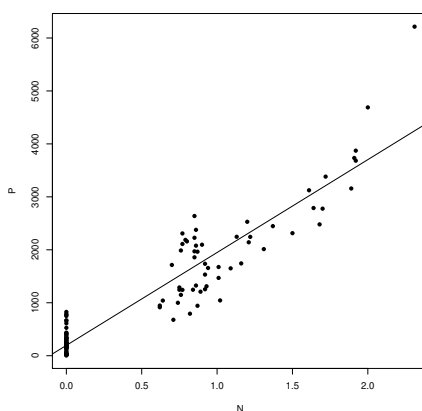
Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	194.93	49.52	3.936	0.000142
N	1753.94	61.07	28.722	< 2e-16

Residual standard error: 408.9 on 115 degrees of freedom
 Multiple R-squared: 0.8777, Adjusted R-squared: 0.8766
 F-statistic: 825 on 1 and 115 DF, p-value: < 2.2e-16

- ii. Eis os comandos e o gráfico resultante:

```
> plot(P ~ N , data=TN025, pch=16)
> abline(TNP.lm)
```

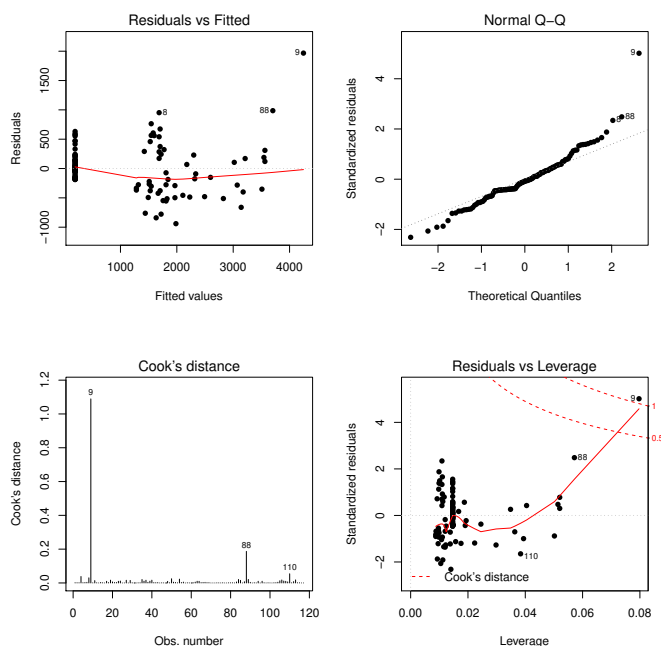


É evidente a existência duma ‘coluna’ de pontos no canto inferior esquerdo do gráfico. Pela escala do eixo horizontal percebe-se que se trata de observações sem a presença do azoto (N). Eis o número de observações nessas condições:

```
> sum(TN025$N==0)
[1] 61
```

Nota: O comando `TN025$N==0` produz um *vector de valores lógicos* TRUE/FALSE indicando quais os elementos da coluna `TN025$N` que verificam a condição “*ser igual a zero*”. O comando `sum` pode agir sobre vectores de valores lógicos, uma vez que converte cada TRUE em “1” e cada FALSE em “0”. Somando, tem-se o número de uns, ou seja, de bacias sem azoto. Havendo 61 bacias hidrográficas (mais de metade) sem azoto, apenas 56 observações ajudam a definir a forma linear da relação.

iii. Eis os quatro gráficos de resíduos e diagnósticos referentes a estes dados:



Os gráficos não colocam em causa os pressupostos do modelo, nem suscitam particulares reparos, com excepção da observação 9, que se destaca em todos. Em particular, a sua distância de Cook é muito elevada: cerca de 1.1, muito acima do limiar 0.5.

Para perceber o que tem de especial essa observação, vamos ver os seus teores de fósforo e azoto, bem como os indicadores de resumo que os permitem avaliar:

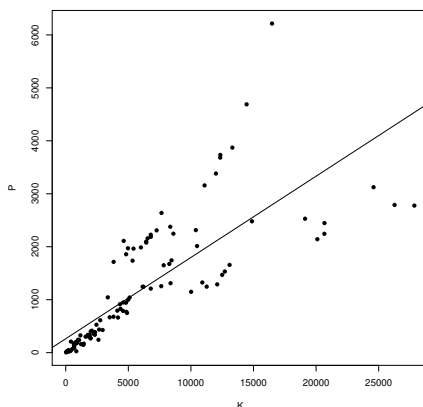
```
> TN025[9,c("P", "N")]
      P      N
9 6214 2.31

> summary(TN025[,c("P", "N")])
      P      N
Min.   : 8    Min.   :0.0000
1st Qu.:182   1st Qu.:0.0000
Median :749   Median :0.0000
Mean   :1114  Mean   :0.5238
3rd Qu.:1964 3rd Qu.:0.8900
Max.   :6214  Max.   :2.3100
```

A observação 9 corresponde à bacia hidrográfica com as maiores concentrações de azoto e fósforo. Esses valores são bastante superiores à da grande maioria das observações, facto confirmado pelo gráfico da subalínea anterior, onde a observação 9 aparece isolada no canto superior direito. O respectivo valor de Cook é muito elevado: trata-se duma observação muito influente, no sentido que a sua exclusão iria alterar bastante a recta ajustada (neste caso, prevê-se que o declive dessa nova recta viesse a ser menor).

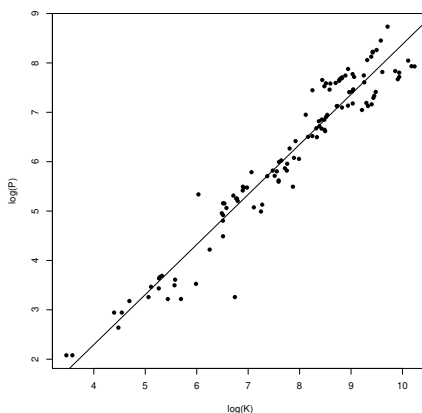
- (c) Substituindo o preditor azoto pelo preditor potássio (K), verifica-se uma espécie de ‘bifurcação’ no gráfico, cuja origem interessaria explorar. Inevitavelmente, o coeficiente de determinação desta nova regressão tem de ser modesto (é $R^2=0.633$) e a recta de regressão resultante será um compromisso entre os dois ‘ramos’ observáveis no gráfico. Já não existe uma ‘coluna’ de pontos, como no gráfico da alínea anterior. Eis os comandos R e o gráfico:

```
> plot(P ~ K , data=TN025, pch=16)
> TNK.lm <- lm(P ~ K , data=TN025)
> abline(TNK.lm)
> summary(TNK.lm)$r.sq
[1] 0.6331735
```



- (d) É pedida agora uma regressão linear simples de log-fósforo sobre log-potássio.

- i. Eis a nuvem de pontos pedida (já com a recta ajustada, só calculada na alínea seguinte):
- ```
> plot(log(P) ~ log(K) , data=TN025, pch=16)
```



O aspecto mais saliente é que a ‘bifurcação’ visível na escala original deixa praticamente de ser visível quando ambas as variáveis são logaritmizadas. A relação linear entre  $\ln(P)$  e  $\ln(K)$  é clara e forte. Tudo indica que a regressão resultante seja de boa qualidade.

ii. Eis a regressão:

```
> TNlnK.lm <- lm(log(P) ~ log(K) , data=TN025)
> summary(TNlnK.lm)
[...]
```

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t ) |
|-------------|----------|------------|---------|----------|
| (Intercept) | -1.7604  | 0.2023     | -8.701  | 2.72e-14 |
| log(K)      | 1.0137   | 0.0254     | 39.906  | < 2e-16  |

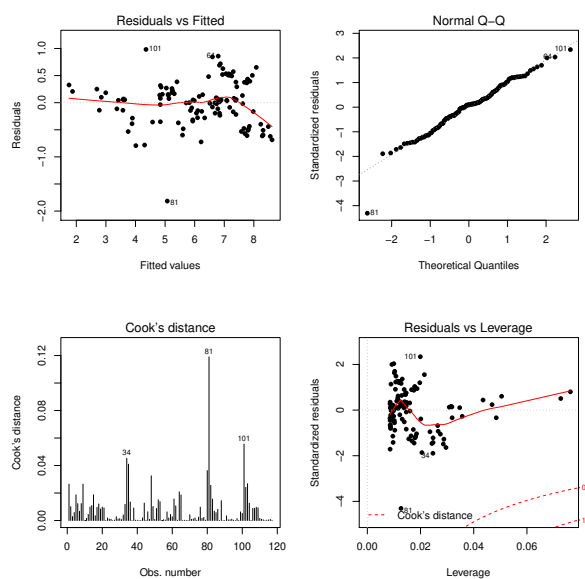
---  
Residual standard error: 0.4242 on 115 degrees of freedom

Multiple R-squared: 0.9326, Adjusted R-squared: 0.9321

F-statistic: 1592 on 1 and 115 DF, p-value: < 2.2e-16

O coeficiente de determinação é elevado: a regressão explica mais de 93% da variabilidade observada nos *log-teores* de fósforo. A transformação logarítmica da variável resposta não permite comparar directamente este valor com o valor obtido na alínea anterior: correspondem a percentagens de coisas diferentes.

iii. Eis os quatro gráficos usuais de resíduos e diagnósticos:

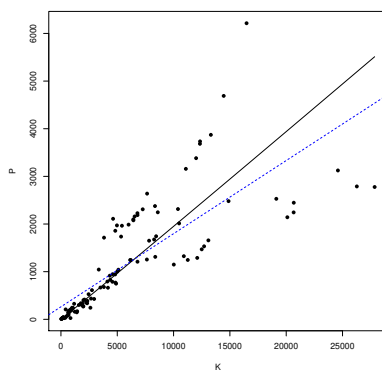


Não há indícios de violação dos pressupostos do modelo. Nenhuma observação tem distância de Cook ou efeito alavanca excessivo. A observação 81 destaca-se nos quatro gráficos, mas sem valores problemáticos. Os valores dos seus log-teores permitem identificá-la como a observação que mais se distancia da recta, sensivelmente a meio da parte inferior no gráfico original. O facto de não ter um valor extremo de log-K reduz a influência desta observação.

```
> log(TN025[81, c("P", "K")])
 P K
81 3.258097 6.741701
```

iv. Uma relação linear entre duas variáveis logaritmizadas corresponde ao modelo potência,  $y = cx^d$ , no nosso caso com  $c = e^{-1.7604} = 0.17198$  e  $d = 1.0137$ . Esta potência muito

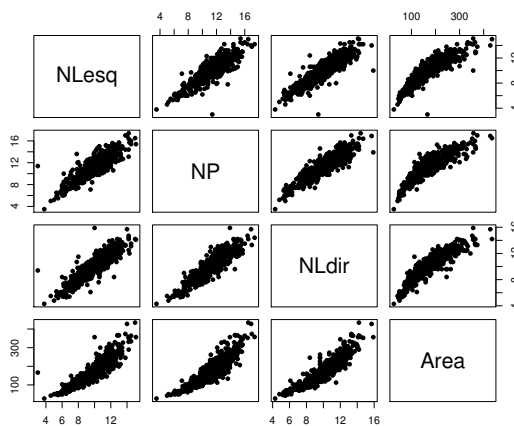
próxima de 1 significa que o gráfico da curva será aproximadamente linear. Eis o gráfico resultante, na escala original, com a curva  $y = 0.17198x^{1.0137}$  (a sólido) e a recta de regressão ajustada na alínea 17c) (a tracejado). Ao regressar para a escala original a ‘bifurcação’ reaparece. A curva potência agora ajustada (que se confirma quase linear, na gama de valores indicados) é assim necessariamente um compromisso entre os dois ‘ramos’, tal como o era a recta de regressão de P sobre K. A curva potência acompanha melhor a tendência dos pontos correspondentes às bacias com valores baixos de P e K.



18. Na *data frame* `videiras`, a primeira coluna indica a casta, pelo que não será de utilidade neste exercício.

(a) Eis o comando para construir as nuvens de pontos pedidas, e o seu resultado:

```
> plot(videiras[,-1], pch=16)
```



Existe uma forte relação linear entre qualquer par de variáveis, pelo que uma regressão linear múltipla de área foliar sobre vários preditores deve ter um coeficiente de determinação elevado. No entanto, nos gráficos que envolvem a variável área, existe alguma evidência de uma ligeira curvatura nas relações com cada comprimento de nervura individual.

(b) Tem-se:

```
> cor(videiras[,-1])
 NLesq NP NLdir Area
NLesq 1.000000 0.878588 0.8870132 0.8902402
```

```

NP 0.878588 1.000000 0.8993985 0.8945700
NLdir 0.8870132 0.8993985 1.0000000 0.8993676
Area 0.8902402 0.8945700 0.8993676 1.0000000

```

Os valores das correlações entre pares de variáveis são todos positivos e bastante elevados, o que confirma as fortes relações lineares evidenciadas nos gráficos.

- (c) Existem  $n$  observações  $\{(x_{1(i)}, x_{2(i)}, x_{3(i)}, Y_i)\}_{i=1}^n$  nas quatro variáveis: a variável resposta área foliar (**Area**, variável aleatória  $Y$ ) e as três variáveis preditoras, associadas aos comprimentos de três nervuras da folha - a principal (variável **NP**,  $X_1$ ), a lateral esquerda (variável **NLesq**,  $X_2$ ) e a lateral direita (variável **NLdir**,  $X_3$ ). Para essas  $n$  observações admite-se que:

- A relação de fundo entre  $Y$  e os três preditores é linear, com variabilidade adicional dada por uma parcela aditiva  $\epsilon_i$  chamada erro aleatório:  
 $Y_i = \beta_0 + \beta_1 x_{1(i)} + \beta_2 x_{2(i)} + \beta_3 x_{3(i)} + \epsilon_i$ , para qualquer  $i = 1, 2, \dots, n$ ;
- os erros aleatórios têm distribuição Normal, de média zero e variância constante:  
 $\epsilon_i \cap \mathcal{N}(0, \sigma^2), \forall i$ ;
- Os erros aleatórios  $\{\epsilon_i\}_{i=1}^n$  são variáveis aleatórias independentes.

- (d) O comando do R que efectua o ajustamento pedido é o seguinte:

```

> videiras.lm <- lm(Area ~ NP + NLesq + NLdir, data=videiras)
> summary(videiras.lm)
[...]
```

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )    |
|-------------|----------|------------|---------|-------------|
| (Intercept) | -168.111 | 5.619      | -29.919 | < 2e-16 *** |
| NP          | 9.987    | 1.192      | 8.380   | 3.8e-16 *** |
| NLesq       | 11.078   | 1.256      | 8.817   | < 2e-16 *** |
| NLdir       | 11.895   | 1.370      | 8.683   | < 2e-16 *** |

---  
Residual standard error: 24.76 on 596 degrees of freedom  
Multiple R-squared: 0.8649, Adjusted R-squared: 0.8642  
F-statistic: 1272 on 3 and 596 DF, p-value: < 2.2e-16

A equação do hiperplano ajustado é assim

$$Area = -168.111 + 9.987 NP + 11.078 NLesq + 11.895 NLdir$$

O valor do coeficiente de determinação é bastante elevado: cerca de 86,49% da variabilidade total nas áreas foliares é explicada por esta regressão linear sobre os comprimentos das três nervuras. Nenhum dos preditores é dispensável sem perda significativa da qualidade do modelo, uma vez que o valor de prova ( $p$ -value) associado aos três testes de hipóteses  $H_0 : \beta_j = 0$  vs.  $H_1 : \beta_j \neq 0$  ( $j = 1, 2, 3$ ) são todos muito pequenos.

O teste de ajustamento global do modelo pode ser formulado assim:

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

**Estatística do teste:**  $F = \frac{QMR}{QMRE} = \frac{n-(p+1)}{p} \frac{R^2}{1-R^2} \cap F_{(p, n-(p+1))}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica (Unilateral direita):** Rej.  $H_0$  se  $F_{calc} > f_{\alpha(p, n-(p+1))} = f_{0.05(3, 596)} \approx 2.62$ .

**Conclusões:** O valor calculado da estatística é dado na listagem produzida pelo R ( $F_{calc} = 1272$ ). Logo, rejeita-se (de forma muito clara) a hipótese nula, que corresponde à hipótese dum modelo inútil. Esta conclusão também resulta directamente da análise do valor de prova ( $p$ -value) associado à estatística de teste calculada:  $p < 2.2 \times 10^{-16}$



corresponde a uma rejeição para qualquer nível de significância usual. Esta conclusão é coerente com o valor bastante elevado de  $R^2$ .

- (e) São pedidos testes envolvendo a hipótese  $\beta_1 = 7$  (não sendo especificada a outra hipótese, deduz-se que seja o complementar  $\beta_1 \neq 7$ ). A hipótese  $\beta_1 = 7$  é uma hipótese simples (um único valor do parâmetro  $\beta_1$ ), que terá de ser colocada na hipótese nula e à qual corresponderá um teste bilateral.

**Hipóteses:**  $H_0 : \beta_1 = 7$  vs.  $H_1 : \beta_1 \neq 7$

**Estatística do Teste:**  $T = \frac{\hat{\beta}_1 - 7}{\hat{\sigma}_{\hat{\beta}_1}} \cap t_{(n-(p+1))}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.01$ .

**Região Crítica:** (Bilateral) Rejeitar  $H_0$  se  $|T_{\text{calc}}| > t_{0.005(596)} \approx 2.584$ .

**Conclusões:** Tem-se  $T_{\text{calc}} = \frac{b_1 - 0}{\hat{\sigma}_{\hat{\beta}_1}} = \frac{9.987 - 7}{1.192} = 2.506 < 2.584$ . Assim, não se rejeita a hipótese nula (que tem o benefício da dúvida), ao nível de significância de 0.01.

Se repetirmos o teste, mas agora utilizando um nível de significância  $\alpha = 0.05$ , apenas a fronteira da região crítica virá diferente. Agora, a regra de rejeição será: rejeitar  $H_0$  se  $|T_{\text{calc}}| > t_{0.025(596)} \approx 1.9640$ . O valor da estatística de teste não se altera ( $T_{\text{calc}} = 2.506$ ), mas este valor pertence agora à região crítica, pelo que ao nível de significância  $\alpha = 0.05$  rejeitamos a hipótese formulada, optando antes por  $H_1 : \beta_1 \neq 7$ . Este exercício ilustra a importância de especificar sempre o nível de significância associado às conclusões do teste.

- (f) É pedido um teste à igualdade de dois coeficientes do modelo, concretamente  $\beta_2 = \beta_3 \Leftrightarrow \beta_2 - \beta_3 = 0$ . Trata-se dum teste à diferença de dois parâmetros, que como foi visto nas aulas, é um caso particular dum teste a uma combinação linear dos parâmetros do modelo. Mais em pormenor, tem-se:

**Hipóteses:**  $H_0 : \beta_2 - \beta_3 = 0$  vs.  $H_1 : \beta_2 - \beta_3 \neq 0$

**Estatística do Teste:**  $T = \frac{(\hat{\beta}_2 - \hat{\beta}_3) - 0}{\hat{\sigma}_{\hat{\beta}_2 - \hat{\beta}_3}} \cap t_{(n-(p+1))}$ , sob  $H_0$

**Nível de significância:**  $\alpha = 0.05$

**Região Crítica:** (Bilateral) Rejeitar  $H_0$  se  $|T_{\text{calc}}| > t_{\alpha/2} (n-(p+1))$

**Conclusões:** Conhecem-se as estimativas  $b_2 = 11.078$  e  $b_3 = 11.895$ , mas precisamos ainda de conhecer o valor do erro padrão associado à estimação de  $\beta_2 - \beta_3$  que, como foi visto nas aulas, é dado por  $\hat{\sigma}_{\hat{\beta}_2 - \hat{\beta}_3} = \sqrt{\hat{V}[\hat{\beta}_2 - \hat{\beta}_3]} = \sqrt{\hat{V}[\hat{\beta}_2] + \hat{V}[\hat{\beta}_3] - 2\widehat{Cov}[\hat{\beta}_2, \hat{\beta}_3]}$ . Assim, precisamos de conhecer as variâncias estimadas de  $\hat{\beta}_2$  e  $\hat{\beta}_3$ , bem como a covariância estimada  $\widehat{cov}[\hat{\beta}_2, \hat{\beta}_3]$ , valores estes que surgem na matriz de (co)variâncias do estimador  $\vec{\beta}$ , que é estimada por  $\hat{V}[\vec{\beta}] = QMRE(\mathbf{X}^t \mathbf{X})^{-1}$ . Esta matriz pode ser calculada no R da seguinte forma:

```
> vcov(videiras.lm)
 (Intercept) NP NLesq NLdir
(Intercept) 31.5707574 -1.0141321 -1.0164689 -0.9051648
NP -1.0141321 1.4200928 -0.6014279 -0.8880395
NLesq -1.0164689 -0.6014279 1.5784886 -0.7969373
NLdir -0.9051648 -0.8880395 -0.7969373 1.8764582
```

Assim,

$$\begin{aligned}\hat{\sigma}_{\hat{\beta}_2 - \hat{\beta}_3} &= \sqrt{\hat{V}[\hat{\beta}_2] + \hat{V}[\hat{\beta}_3] - 2\widehat{Cov}[\hat{\beta}_2, \hat{\beta}_3]} \\ &= \sqrt{1.5784886 + 1.8764582 - 2 \times (-0.7969373)} = \sqrt{5.048821} = 2.246958,\end{aligned}$$

pelo que  $T_{\text{calc}} = \frac{11.078 - 11.895}{2.246958} = -0.3636027$ . Como  $|T_{\text{calc}}| < t_{0.025(596)} \approx 1.9640$ , não se rejeita  $H_0$  ao nível de significância de 0.05, isto é, admite-se que  $\beta_2 = \beta_3$ . No contexto do problema, não se rejeitou a hipótese que a variação média provocada na área foliar seja igual, quer se aumente a nervura lateral esquerda ou a nervura lateral direita em 1cm (mantendo as restantes nervuras de igual comprimento).

- (g) i. Substituindo na equação do hiperplano ajustado, obtido na alínea 18d, obtêm-se os seguintes valores estimados:

- *Folha 1*:  $\widehat{Area} = -168.111 + 9.987 \times 12.1 + 11.078 \times 11.6 + 11.895 \times 11.9 = 222.787 \text{ cm}^2$ ;
- *Folha 2*:  $\widehat{Area} = -168.111 + 9.987 \times 10.6 + 11.078 \times 10.1 + 11.895 \times 9.9 = 167.3995 \text{ cm}^2$ ;
- *Folha 3*:  $\widehat{Area} = -168.111 + 9.987 \times 15.1 + 11.078 \times 14.9 + 11.895 \times 14.0 = 314.2849 \text{ cm}^2$ ;

Com recurso ao comando `predict` do R, estas três áreas ajustadas obtêm-se da seguinte forma:

```
> predict(videiras.lm, new=data.frame(NP=c(12.1,10.6,15.1), NLesq=c(11.6,10.1,14.9),
+ NLdir=c(11.9, 9.9, 14.0)))
 1 2 3
222.7762 167.3903 314.2715
```

Novamente, algumas pequenas discrepâncias nas casas decimais finais resultam de erros de arredondamento.

- ii. Estes intervalos de confiança para  $\mu_{Y|X} = E[Y|X_1 = x_1, X_2 = x_2, X_3 = x_3]$  (com os valores de  $x_1$ ,  $x_2$  e  $x_3$  indicados no enunciado, para cada uma das três folhas) obtêm-se subtraindo e somando aos valores ajustados obtidos na subalínea anterior a semi-amplitude do IC, dada por  $t_{\alpha/2(n-(p+1))} \cdot \hat{\sigma}_{\hat{\mu}_{Y|X}}$ , sendo  $\hat{\sigma}_{\hat{\mu}_{Y|X}} = \sqrt{QMRE \cdot \mathbf{a}^t(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{a}}$  onde os vectores  $\mathbf{a}$  são os vectores da forma  $\mathbf{a} = (1, x_1, x_2, x_3)$ . Estas contas, algo trabalhosas, resultam fáceis recorrendo de novo ao comando `predict` do R, mas desta vez com o argumento `int="conf"`, como indicado de seguida:

```
> predict(videiras.lm, new=data.frame(NP=c(12.1,10.6,15.1), NLesq=c(11.6,10.1,14.9),
+ NLdir=c(11.9, 9.9, 14.0)), int="conf")
 fit lwr upr
1 222.7762 219.1776 226.3747
2 167.3903 164.9215 169.8590
3 314.2715 308.4607 320.0823
```

Assim, tem-se para cada folha, os seguintes intervalos a 95% de confiança para  $\mu_{Y|X}$ :

- *Folha 1*: ] 219.1776 , 226.3747 [;
- *Folha 2*: ] 164.9215 , 169.8590 [;
- *Folha 3*: ] 308.4607 , 320.0823 [.

Repare-se como a amplitude de cada intervalo é diferente, uma vez que depende de informação específica para cada folha (dada pelo vector  $\mathbf{a}$  dos valores dos preditores).

- iii. Sabemos que os intervalos de predição têm uma forma análoga aos intervalos de confiança para  $E[Y|X]$ , mas com uma maior amplitude, associada à variabilidade adicional de observações individuais, a que corresponde  $\hat{\sigma}_{\text{indiv}} = \sqrt{QMRE \cdot [1 + \mathbf{a}^t(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{a}]}$ . De novo, recorreremos ao comando `predict`, desta vez com o argumento `int="pred"`:

```

> predict(videiras.lm, new=data.frame(NP=c(12.1,10.6,15.1),NLesq=c(11.6,10.1,14.9),
+ Nmdir=c(11.9, 9.9, 14.0)), int="pred")
 fit lwr upr
1 222.7762 174.0206 271.5318
2 167.3903 118.7050 216.0755
3 314.2715 265.3029 363.2401

```

Assim, têm-se os seguintes intervalos de predição a 95% para os três valores de  $Y$ :

- *Folha 1*: ] 174.0206 , 271.5318 [;
- *Folha 2*: ] 118.7050 , 216.0755 [;
- *Folha 3*: ] 265.3029 , 363.2401 [.

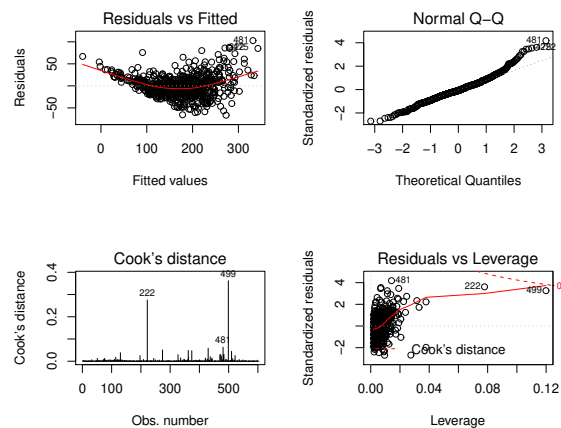
A amplitude maior destes intervalos deve-se ao valor elevado do Quadrado Médio Residual, que estima a variabilidade das observações individuais de  $Y$  em torno da recta.

- (h) Recorremos de novo ao R para construir os gráficos de resíduos. O primeiro dos dois comandos seguintes destina-se a dividir a janela gráfica numa espécie de matriz  $2 \times 2$ :

```

> par(mfrow=c(2,2))
> plot(videiras.lm, which=c(1,2,4,5))

```



O gráfico do canto superior esquerdo é o gráfico dos resíduos usuais ( $e_i$ ) vs. valores ajustados ( $\hat{y}_i$ ). Neste gráfico são visíveis dois problemas: uma tendência para a curvatura (já detectado nos gráficos da variável resposta contra cada preditor individual), que indica que o modelo linear pode não ser a melhor forma de relacionar área foliar com os comprimentos das nervuras; e uma forma em funil que sugere que a hipótese de homogeneidade das variâncias dos erros aleatórios pode não ser a mais adequada.

No canto superior direito tem-se um *qq-plot*, de quantis empíricos vs. quantis teóricos duma Normal reduzida. A ser verdade a hipótese de Normalidade dos erros aleatórios, seria de esperar uma disposição linear dos pontos neste gráfico. É visível, sobretudo na parte direita do gráfico, um afastamento relativamente forte de muitas observações a esta linearidade, sugerindo problemas com a hipótese de Normalidade.

No canto inferior esquerdo tem-se um diagrama de barras com as distâncias de Cook de cada observação. Embora nenhuma observação ultrapasse o limiar de guarda  $D_i > 0.5$ , duas observações têm um valor considerável da distância de Cook: a observação 499, com  $D_{499}$  próximo de 0.4 e a observação 222, com distância de Cook próxima de 0.3. Estas duas observações merecem especial atenção. Finalmente, o gráfico no canto inferior direito

relaciona resíduos (internamente) estandardizados (eixo vertical) com valor do efeito alavanca (eixo horizontal) e também com as distâncias de Cook (sendo traçadas pelo R linhas de igual distância de Cook, para alguns valores particularmente elevados, como 0.5 ou 1). Este gráfico ilustra que as duas observações com maior distância de Cook (499 e 222) têm valores relativamente elevados, quer dos resíduos estandardizados, quer do efeito alavanca. O efeito alavanca médio, neste ajustamento de  $n = 600$  observações a um modelo com  $p + 1 = 4$  parâmetros é  $\bar{h} = \frac{4}{600} = 0.006667$  e as duas observações referidas têm os maiores efeitos alavanca das  $n = 600$  observações, respectivamente, próximos de 0.12 e 0.08. Já a observação 481, igualmente identificada no gráfico, tem o maior resíduo estandardizado de qualquer observação, mas com um valor relativamente discreto do efeito alavanca, acaba por não ser uma observação influente (como se pode confirmar no gráfico anterior).

Este exemplo confirma que são diferentes os conceitos de resíduo elevado, observação influente e elevado valor do efeito alavanca (*leverage*). Os valores das variáveis observadas, nas folhas 222 e 499, revelam um desequilíbrio nos comprimentos das nervuras laterais, sendo em ambos os casos a nervura lateral direita muito mais comprida do que a esquerda. Além disso, ambas as folhas têm uma das nervuras laterais de comprimento extremo: no caso da folha 222 tem-se a maior nervura lateral direita de qualquer das 600 folhas, enquanto que a folha 499 tem a mais pequena de todas as nervuras laterais esquerdas. Assim, trata-se de folhas com formas irregulares, diferentes da generalidade das folhas analisadas.

Este exercício visa chamar a atenção que *um modelo de regressão com um ajustamento bastante forte pode revelar, no estudo dos resíduos, problemas* que levantam dúvidas sobre a validade das conclusões inferenciais (testes de hipóteses, intervalos de confiança e predição) obtidas nas alíneas anteriores.

- (i) O pedido de logaritmizar previamente as variáveis envolvidas no estudo faz sentido, tendo em conta a curvilinearidade sugerida pelo gráfico de resíduos da alínea anterior (18h). Eis o resultado do ajustamento pedido:

```
> videiraslog.lm <- lm(log(Area) ~ log(NP) + log(NLesq) + log(NLdir), data=videiras)
> summary(videiraslog.lm)
```

```
Call: lm(formula = log(Area) ~ log(NP) + log(NLesq) + log(NLdir), data = videiras)
[...]
```

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 0.40983  | 0.06136    | 6.679   | 5.52e-11 *** |
| log(NP)     | 0.72660  | 0.06574    | 11.052  | < 2e-16 ***  |
| log(NLesq)  | 0.57049  | 0.05649    | 10.100  | < 2e-16 ***  |
| log(NLdir)  | 0.71077  | 0.06780    | 10.484  | < 2e-16 ***  |

---

Residual standard error: 0.1259 on 596 degrees of freedom

Multiple R-squared: 0.9081, Adjusted R-squared: 0.9076

F-statistic: 1963 on 3 and 596 DF, p-value: < 2.2e-16

Não é legítimo procurar comparar directamente o coeficiente de determinação deste modelo,  $R^2 = 0.9081$ , e o coeficiente de determinação do modelo análogo sem a logaritmização (alínea 18d),  $R^2 = 0.8649$ , uma vez que a escala onde são medidos os resíduos são diferentes, nos dois casos. Apenas podemos afirmar que o modelo agora ajustado explica mais de 90% da variância dos valores observados *das log-áreas* foliares. A equação do hiperplano ajustado é da forma  $\ln(y) = b_0 + b_1 \ln(x_1) + b_2 \ln(x_2) + b_3 \ln(x_3)$ , sendo  $y$  a *Area*,  $x_1$  a variável NP,  $x_2$  a variável NLesq, e  $x_3$  a variável NLdir, e tendo  $b_0 = 0.40983$ ,  $b_1 = 0.72660$ ,  $b_2 = 0.57049$

e  $b_3=0.71077$ . Em termos das variáveis originais esta relação corresponde a:

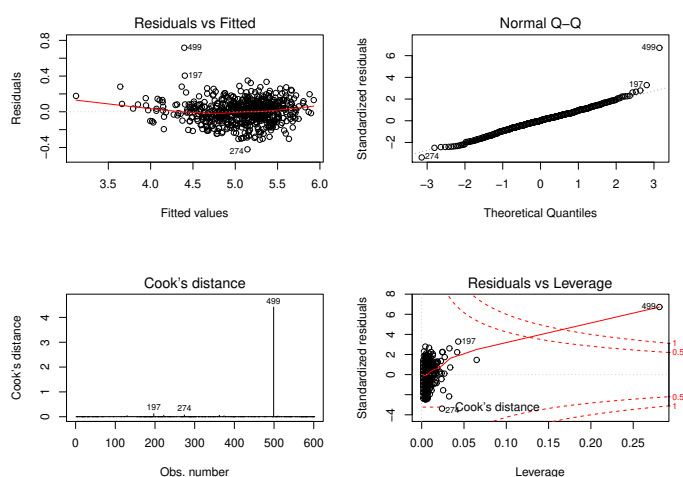
$$\begin{aligned} \ln(y) = b_0 + b_1 \ln(x_1) + b_2 \ln(x_2) + b_3 \ln(x_3) &\Leftrightarrow y = \exp^{b_0+b_1 \ln(x_1)+b_2 \ln(x_2)+b_3 \ln(x_3)} \\ &\Leftrightarrow y = e^{b_0} e^{b_1 \ln(x_1)} e^{b_2 \ln(x_2)} e^{b_3 \ln(x_3)} \\ &\Leftrightarrow y = e^{b_0} e^{\ln(x_1^{b_1})} e^{\ln(x_2^{b_2})} e^{\ln(x_3^{b_3})} \\ &\Leftrightarrow y = e^{b_0} x_1^{b_1} x_2^{b_2} x_3^{b_3} \end{aligned}$$

Logo o modelo ajustado tem a seguinte equação:

$$Area = 1.506562 NP^{0.72660} NLesq^{0.57049} NLdir^{0.71077} .$$

(j) Proceda-se como na alínea 18h) e obtém-se:

```
> par(mfrow=c(2,2))
> plot(videiraslog.lm, which=c(1,2,4,5))
```



Os problemas identificados na alínea 18h) foram em boa medida corrigidos. Assim, a logaritização das quatro variáveis revelou ser uma opção adequada. O gráfico de resíduos usuais  $e_i$  contra valores ajustados  $\hat{y}_i$  mostra agora uma dispersão dos pontos numa banda horizontal em torno do valor médio zero, tendo praticamente desaparecido, quer a curvatura, quer a forma em funil. Assim, a linearidade da relação entre as variáveis logaritimizadas, bem como a homogeneidade das variâncias dos respectivos erros aleatórios são pressupostos admissíveis. Da mesma forma, o *qq-plot* do canto superior direito mostra que (à excepção da observação 499) tem-se uma boa linearidade, sustentando o pressuposto de Normalidade dos erros aleatórios. No canto inferior esquerdo, a observação 499 surge de novo destacada, com uma enormíssima distância de Cook, superior a 4, e portanto muito superior ao limiar de guarda 0.5. Assim, esta observação tem uma enorme influência na regressão ajustada, e a sua exclusão provocaria alterações importantes, quer nos coeficientes ajustados  $b_j$ , quer nos valores resultantes de  $\hat{y}_i$ . Essa mesma indicação é dada no quarto e último gráfico, onde (graças ao elevadíssimo valor de  $D_{499}$ ) são visíveis dois pares de isolinhas de Cook, correspondentes aos limiares 0.5 e 1. Registe-se ainda, nos dois gráficos da direita, como a observação 499 tem um enorme resíduo (internamente) estandardizado, com  $R_{499} > 6$ , bem como um efeito alavanca razoavelmente elevado (que nenhuma outra observação acompanha). A discordante observação 499 (que é, simultaneamente uma observação atípica, influente e de valor razoavelmente elevado do efeito alavanca) já foi discutida anteriormente. Tratando-se de uma folha com uma muito evidente assimetria (possivelmente correspondente a um

erro de medição/registo, ou então danificada por alguma razão), haverá espaço para discutir a sua eventual exclusão do modelo, podendo argumentar-se que o modelo destina-se a ser usado com folhas de videira não danificadas ou excessivamente irregulares. A título de curiosidade, registe-se o resultado de reajustar o modelo, apenas com as 599 folhas restantes:

```
> summary(lm(log(Area) ~ log(NP) + log(NLesq) + log(NLdir), data=videiras[-499,]))

Call: lm(formula = log(Area) ~ log(NP) + log(NLesq) + log(NLdir), data = videiras[-499,])
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.38758 0.05912 6.555 1.2e-10 ***
log(NP) 0.60695 0.06553 9.262 < 2e-16 ***
log(NLesq) 0.80654 0.06399 12.604 < 2e-16 ***
log(NLdir) 0.60978 0.06681 9.127 < 2e-16 ***

Residual standard error: 0.1211 on 595 degrees of freedom
Multiple R-squared: 0.9151, Adjusted R-squared: 0.9146
F-statistic: 2137 on 3 and 595 DF, p-value: < 2.2e-16
```

Repare-se na alteração substancial dos valores estimados dos quatro parâmetros, e em especial dos coeficientes dos log-comprimentos das nervuras, uma alteração que confirma que a observação 499 era muito influente.

19. (a) Eis a regressão linear múltipla de rendimento de milho sobre todos os preditores:

```
> summary(lm(y ~ . , data=milho))
[...]
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) 51.03036 85.73770 0.595 0.557527
x1 0.87691 0.18746 4.678 0.000104 ***
x2 0.78678 0.43036 1.828 0.080522 .
x3 -0.46017 0.42906 -1.073 0.294617
x4 -0.77605 1.05512 -0.736 0.469464
x5 0.48279 0.57352 0.842 0.408563
x6 2.56395 1.38032 1.858 0.076089 .
x7 0.05967 0.71881 0.083 0.934556
x8 0.40590 1.03322 0.393 0.698045
x9 -0.65951 0.67034 -0.984 0.335426

Residual standard error: 7.815 on 23 degrees of freedom
Multiple R-squared: 0.7476, Adjusted R-squared: 0.6488
F-statistic: 7.569 on 9 and 23 DF, p-value: 4.349e-05
```

Não sendo um ajustamento excelente, as variáveis preditivas conseguem explicar quase 75% da variabilidade nos rendimentos. Um teste de ajustamento global rejeita a hipótese nula (inutilidade do modelo) com um valor de prova  $p=0.00004349$ .

- (b) O coeficiente de determinação modificado é dado na listagem produzida pelo R:  $R_{mod}^2 = 0.6488$ . Define-se como  $R_{mod}^2 = 1 - \frac{QMRE}{QMT} = 1 - \frac{SQRE}{SQT} \cdot \frac{n-1}{n-(p+1)} = 1 - (1 - R^2) \cdot \frac{n-1}{n-(p+1)}$ . A diferença óbvia, neste exercício, entre os valores do  $R^2$  usual e de  $R_{mod}^2$  resulta de termos um valor de  $R^2$  (usual) modesto e ajustado com um número de observações ( $n=33$ ) pequena face ao número de parâmetros do modelo ( $p+1=10$ ). Pela última das expressões acima para  $R_{mod}^2$ , vemos que o factor multiplicativo  $\frac{n-1}{n-(p+1)} = \frac{32}{23} = 1.3913$ . Assim, a distância do  $R^2$  usual em relação ao seu máximo ( $1 - R^2 = 0.2524$ ) é aumentado em cerca de 40% antes de ser subtraído de novo ao valor máximo 1, pelo que  $R_{mod}^2 = 1 - 0.2524 \times 1.3913 = 1 - 0.3512 = 0.6488$ . Em geral,  $R_{mod}^2$  penaliza modelos ajustados com relativamente poucas observações (em

relação ao número de parâmetros do modelo), em especial quando o valor de  $R^2$  não é muito elevado. Por outras palavras,  $R_{mod}^2$  penaliza modelos com ajustamentos modestos, baseados em relativamente pouca informação, face à complexidade do modelo.

(c) Eis o resultado do ajustamento pedido, sem o preditor  $x_1$ :

```
> summary(lm(y ~ . - x1 , data=milho))
[...]
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) 192.387333 109.724668 1.753 0.0923 .
x2 0.305508 0.571461 0.535 0.5978
x3 -0.469256 0.586748 -0.800 0.4317
x4 -1.526474 1.426129 -1.070 0.2951
x5 -0.133203 0.763345 -0.174 0.8629
x6 3.312695 1.874882 1.767 0.0900 .
x7 -1.580293 0.858146 -1.842 0.0779 .
x8 1.239484 1.391780 0.891 0.3820
x9 -0.008387 0.896726 -0.009 0.9926

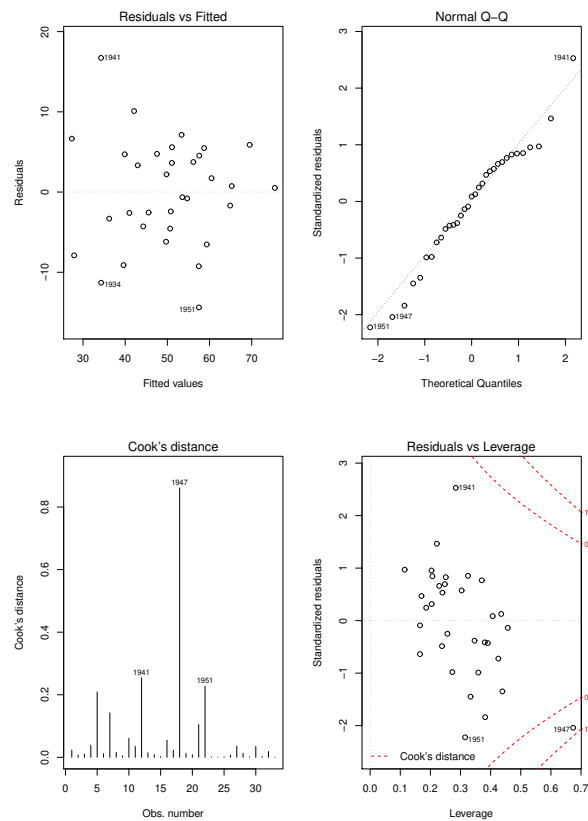
Residual standard error: 10.69 on 24 degrees of freedom
Multiple R-squared: 0.5074, Adjusted R-squared: 0.3432
F-statistic: 3.091 on 8 and 24 DF, p-value: 0.01524
```

A exclusão do preditor  $x_1$  provocou uma quebra acentuada no valor do coeficiente de determinação, que é agora apenas  $R^2=0.5074$  (repare-se como  $R_{mod}^2=0.3432$  ainda se distancia mais do  $R^2$  usual). Assim, este modelo sem o preditor  $x_1$  apenas explica cerca de metade da variabilidade nos rendimentos. Outro facto saliente é a grande perturbação nos valores ajustados dos parâmetros (quando comparados com o modelo com todos os preditores).

Este enorme impacto da exclusão do preditor  $x_1$  é digno de nota, porque essa variável preditora é apenas um contador dos anos que passam. Há dois aspectos a salientar:

- o preditor  $x_1$  funciona aqui como uma variável substituta (*proxy variable*, em inglês) para um grande número de outras variáveis, muitas das quais de difícil quantificação, tais como desenvolvimentos técnicos ou tecnológicos associados à cultura do milho nos anos em questão. Revela ser um indicador simples para levar em conta aspectos não meteorológicos que, nos anos em questão, tiveram grande impacto na produção e que não eram contemplados pelos restantes preditores.
- este exemplo ilustra como os modelos estudam *associações estatísticas*, o que não é sinónimo com *relações de causa e efeito*. No ajustamento do modelo com todos os preditores, a estimativa do coeficiente da variável  $x_1$  é  $b_1 = 0.87691$ . Tendo em conta a natureza e unidades de medida das variáveis, podemos afirmar que, a cada ano que passa (e mantendo as restantes variáveis constantes) o valor da produção aumenta, em média,  $0.87691$  *bushels/acre*. Mas não faz evidentemente sentido dizer que cada ano que passa *provoca* esse aumento na produção. Não é a mera passagem do tempo que *causa* a produção. Pode existir uma relação de causa e efeito entre alguns preditores e a variável resposta, mas pode apenas existir *associação*, como neste caso. A existência, ou não, de uma relação de causa e efeito nunca poderá ser afirmada pela via estatística, mas apenas com base nos conhecimentos teóricos associados aos fenómenos sob estudo.

Quanto ao estudo dos resíduos, eis os gráficos produzidos com as opções 1, 2, 4 e 5 do comando `plot` do R:



O gráfico de resíduos usuais *vs.* valores ajustados  $\hat{y}_i$  (no canto superior esquerdo) não apresenta qualquer padrão digno de registo, dispersando-se os resíduos numa banda horizontal. Assim, nada sugere que não se verifiquem os pressupostos de linearidade e de homogeneidade de variâncias, admitidos no modelo RLM. Analogamente, no *qq-plot* comparando quantis teóricos duma Normal reduzida e quantis empíricos (canto superior direito), existe linearidade aproximada dos pontos, pelo que a hipótese de Normalidade dos erros aleatórios também parece admissível. Já no diagrama de barras das distâncias de Cook (canto inferior esquerdo) há um facto digno de registo: a observação correspondente ao ano 1947 tem um valor elevadíssimo da distância de Cook (superior a 0.8), pelo que se trata dum ano muito influente no ajustamento do modelo. Dado o elevado número de variáveis preditoras, não é possível visualizar a nuvem de pontos associada aos dados, mas uma análise mais atenta da tabela de valores observados (disponível no enunciado) sugere possíveis causas para este facto. O ano de 1947 teve uma precipitação pré-Junho particularmente intensa, a que se seguiu um mês de Agosto anormalmente quente e seco (nas três variáveis registam-se observações extremas, para os anos observados). O valor muito elevado da distância de Cook indica que a exclusão deste ano do conjunto de dados provocaria alterações importantes no modelo ajustado. Finalmente, o gráfico de resíduos internamente estandardizados ( $R_i$ ) *vs.* valores do efeito alavanca ( $h_{ii}$ ) confirmam a elevada distância de Cook da observação correspondente a 1947, e mostram que ela resulta dum resíduo internamente estandardizado relativamente grande, em valor absoluto (embora não extraordinariamente grande), mas sobretudo dum valor muito elevado (cerca de 0.7) do efeito alavanca. Este último valor sugere que esta observação está a “atrair” o hiperplano ajustado, facto que ajuda a esconder





Apenas aceitando trabalhar com uma probabilidade de cometer o erro de Tipo I maior, por exemplo  $\alpha = 0.10$ , é que seria possível rejeitar  $H_0$  e considerar os modelos como tendo ajustamentos significativamente diferentes.

Esta conclusão sugere a possibilidade de ter, já em finais de Junho, previsões de produção que expliquem quase dois terços da variabilidade observada na produção. No entanto, deve recordar-se que se trata dum modelo ajustado com relativamente poucas observações.

- (e) Vamos aplicar o algoritmo de exclusão sequencial, baseado nos testes  $t$  aos coeficientes  $\beta_j$  e usando um nível de significância  $\alpha = 0.10$ .

Partindo do ajustamento do modelo com todos os preditores, efectuado na alínea 19a), conclui-se que há várias variáveis candidatas a sair (os  $p$ -values correspondentes aos testes a  $\beta_j = 0$  são superiores ao limiar acima indicado). De entre estas, é a variável  $x_7$  que tem de longe o maior  $p$ -value, pelo que é a primeira variável a excluir.

Após a exclusão do preditor  $x_7$  é necessário re-ajustar o modelo:

```
> summary(lm(y ~ . - x7, data=milho))
[...]
```

|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 54.8704  | 70.6804    | 0.776   | 0.4451       |
| x1          | 0.8693   | 0.1602     | 5.425   | 1.42e-05 *** |
| x2          | 0.7751   | 0.3983     | 1.946   | 0.0634 .     |
| x3          | -0.4590  | 0.4199     | -1.093  | 0.2852       |
| x4          | -0.7982  | 0.9995     | -0.799  | 0.4324       |
| x5          | 0.4814   | 0.5613     | 0.858   | 0.3996       |
| x6          | 2.5245   | 1.2687     | 1.990   | 0.0581 .     |
| x8          | 0.4137   | 1.0074     | 0.411   | 0.6849       |
| x9          | -0.6426  | 0.6252     | -1.028  | 0.3143       |

```

Residual standard error: 7.652 on 24 degrees of freedom
Multiple R-squared: 0.7475, Adjusted R-squared: 0.6633
F-statistic: 8.882 on 8 and 24 DF, p-value: 1.38e-05
```

Assinale-se que o valor do coeficiente de determinação quase não se alterou com a exclusão de  $x_7$ . Continuam a existir várias variáveis com valor de prova superiores ao limiar estabelecido, e de entre estas é a variável  $x_8$  que tem o maior  $p$ -value:  $p = 0.6849$ . Exclui-se essa variável e ajusta-se novamente o modelo.

```
> summary(lm(y ~ . - x7 - x8, data=milho))
[...]
```

|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 58.4750  | 68.9575    | 0.848   | 0.4045       |
| x1          | 0.8790   | 0.1558     | 5.641   | 7.17e-06 *** |
| x2          | 0.8300   | 0.3689     | 2.250   | 0.0335 *     |
| x3          | -0.4592  | 0.4128     | -1.112  | 0.2765       |
| x4          | -0.8354  | 0.9787     | -0.854  | 0.4015       |
| x5          | 0.5287   | 0.5401     | 0.979   | 0.3370       |
| x6          | 2.4392   | 1.2306     | 1.982   | 0.0586 .     |
| x9          | -0.7254  | 0.5819     | -1.247  | 0.2240       |

```

Residual standard error: 7.523 on 25 degrees of freedom
```

---

Multiple R-squared: 0.7457, Adjusted R-squared: 0.6745  
F-statistic: 10.47 on 7 and 25 DF, p-value: 4.333e-06

O valor de  $R^2$  mantém-se próximo do original e continuam a existir variáveis candidatas a sair do modelo. De entre estas, é o preditor  $x_4$  que tem o maior  $p$ -value ( $p = 0.4015$ ), pelo que será o próximo preditor a excluir. O re-ajustamento do modelo sem os três preditores já excluídos ( $x_7$ ,  $x_8$  e  $x_4$ ) produz os seguintes resultados:

```
> summary(lm(y ~ . - x7 - x8 - x4, data=milho))
[...]
```

|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 37.9486  | 64.2899    | 0.590   | 0.5601       |
| x1          | 0.8854   | 0.1548     | 5.718   | 5.11e-06 *** |
| x2          | 0.7685   | 0.3599     | 2.135   | 0.0423 *     |
| x3          | -0.3603  | 0.3941     | -0.914  | 0.3690       |
| x5          | 0.6338   | 0.5231     | 1.212   | 0.2366       |
| x6          | 2.7275   | 1.1772     | 2.317   | 0.0286 *     |
| x9          | -0.6829  | 0.5767     | -1.184  | 0.2471       |

```

Residual standard error: 7.484 on 26 degrees of freedom
Multiple R-squared: 0.7383, Adjusted R-squared: 0.6779
F-statistic: 12.23 on 6 and 26 DF, p-value: 1.624e-06
```

Após a exclusão de três preditores, o coeficiente de determinação continua próximo do valor original:  $R^2 = 0.7383$ . Esta quebra pequena reflecte os valores elevados dos  $p$ -values associados aos preditores excluídos. Mas há mais preditores candidatos à exclusão, sendo  $x_3$  a próxima variável a excluir do lote de preditores ( $p=0.3690 > 0.10$ ).

```
> summary(lm(y ~ . - x7 - x8 - x4 - x3, data=milho))
[...]
```

|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 39.3646  | 64.0755    | 0.614   | 0.5441       |
| x1          | 0.8870   | 0.1544     | 5.747   | 4.13e-06 *** |
| x2          | 0.7562   | 0.3586     | 2.109   | 0.0444 *     |
| x5          | 0.4725   | 0.4910     | 0.962   | 0.3444       |
| x6          | 2.4893   | 1.1445     | 2.175   | 0.0386 *     |
| x9          | -0.8320  | 0.5515     | -1.509  | 0.1430       |

```

Residual standard error: 7.461 on 27 degrees of freedom
Multiple R-squared: 0.7299, Adjusted R-squared: 0.6799
F-statistic: 14.59 on 5 and 27 DF, p-value: 5.835e-07
```

Há ainda candidatos à exclusão, sendo  $x_5$  a exclusão seguinte.

```
> summary(lm(y ~ . - x7 - x8 - x4 - x3 - x5, data=milho))
[...]
```

|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 87.1589  | 40.4371    | 2.155   | 0.0399 *     |
| x1          | 0.8519   | 0.1498     | 5.688   | 4.25e-06 *** |
| x2          | 0.5989   | 0.3187     | 1.879   | 0.0707 .     |

```

x6 2.3613 1.1353 2.080 0.0468 *
x9 -0.9755 0.5302 -1.840 0.0764 .

Residual standard error: 7.451 on 28 degrees of freedom
Multiple R-squared: 0.7206, Adjusted R-squared: 0.6807
F-statistic: 18.06 on 4 and 28 DF, p-value: 1.954e-07

```

Tendo em conta que fixámos o limiar de exclusão no nível de significância  $\alpha = 0.10$ , não há mais variáveis candidatas à exclusão, pelo que o algoritmo termina aqui. O modelo final escolhido pelo algoritmo tem quatro preditores ( $x_1$ ,  $x_2$ ,  $x_6$  e  $x_9$ ), e um coeficiente de determinação  $R^2 = 0.7206$ . Ou seja, com menos de metade dos preditores iniciais, apenas se perdeu 0.027 no valor de  $R^2$ .

O valor relativamente alto ( $\alpha = 0.10$ ) do nível de significância usado é aconselhável, na aplicação deste algoritmo, uma vez que variáveis cujo *p-value* cai abaixo deste limiar podem, se excluídas, gerar quebras mais pronunciadas no valor de  $R^2$ . Tal facto é ilustrado pela exclusão de  $x_9$  (a exclusão seguinte, caso se tivesse optado por um limiar  $\alpha = 0.05$ ):

```

> summary(lm(y ~ . - x7 - x8 - x4 - x3 - x5 - x9, data=milho))
[...]
Residual standard error: 7.752 on 29 degrees of freedom
Multiple R-squared: 0.6869, Adjusted R-squared: 0.6545
F-statistic: 21.2 on 3 and 29 DF, p-value: 1.806e-07

```

Dado o número de exclusões efectuadas, pode desejar-se fazer um teste  $F$  parcial, comparando o submodelo final produzido pelo algoritmo e o modelo original com todos os preditores:

```

> anova(milhoAlgExc.lm, milho.lm)
Analysis of Variance Table

Model 1: y ~ x1 + x2 + x6 + x9
Model 2: y ~ x1 + x2 + x3 + x4 + x5 + x6 + x7 + x8 + x9
 Res.Df RSS Df Sum of Sq F Pr(>F)
1 28 1554.6
2 23 1404.7 5 149.9 0.4909 0.7796

```

O *p-value* muito elevado ( $p = 0.7796$ ) indica que não se rejeita a hipótese de modelo e submodelo serem equivalentes.

Como foi indicado nas aulas, existe uma função do R, a função `step`, que automatiza um algoritmo de exclusão sequencial, mas utilizando o valor do Critério de Informação de Akaike (AIC) como critério de exclusão dum preditor em cada passo do algoritmo. Esta função produz neste exemplo o mesmo submodelo final, como se pode constatar na parte final desta listagem:

```

> step(milho.lm)
Start: AIC=143.79
y ~ x1 + x2 + x3 + x4 + x5 + x6 + x7 + x8 + x9
[...]
Step: AIC=137.13
y ~ x1 + x2 + x6 + x9
 Df Sum of Sq RSS AIC
<none> 1554.6 137.13

```

```

- x9 1 187.95 1742.6 138.90
- x2 1 196.01 1750.6 139.05
- x6 1 240.20 1794.8 139.87
- x1 1 1796.22 3350.8 160.47
Call: lm(formula = y ~ x1 + x2 + x6 + x9, data = milho)
Coefficients:
(Intercept) x1 x2 x6 x9
 87.1589 0.8519 0.5989 2.3613 -0.9755

```

Refira-se que as variáveis meteorológicas mais associadas à previsão da produção são a precipitação pré-Junho ( $x_2$ ), a precipitação em Julho ( $x_6$ ) e a temperatura em Agosto ( $x_9$ ). Finalmente, refira-se que, caso esteja disponível *software* adequado, pode recorrer-se a uma pesquisa completa de todos os subconjuntos, a fim de escolher os melhores, para cada número  $k$  de preditores. Como referido nas aulas, o módulo `leaps` do R disponibiliza um comando de igual nome para fazer essas escolhas. Eis os comandos e a listagem produzida, para o conjunto de dados deste Exercício.

```

> library(leaps)
> leaps(y=milho$y , x=milho[,-10], method="r2", nbest=1)
$which
 1 2 3 4 5 6 7 8 9
1 TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
2 TRUE FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE
3 TRUE TRUE FALSE FALSE FALSE TRUE FALSE FALSE FALSE
4 TRUE TRUE FALSE FALSE FALSE TRUE FALSE FALSE TRUE
5 TRUE TRUE FALSE FALSE TRUE TRUE FALSE FALSE TRUE
6 TRUE TRUE TRUE FALSE TRUE TRUE FALSE FALSE TRUE
7 TRUE TRUE TRUE TRUE TRUE TRUE FALSE FALSE TRUE
8 TRUE TRUE TRUE TRUE TRUE TRUE FALSE TRUE TRUE
9 TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
[...]
$r2
[1] 0.5633857 0.6566246 0.6868757 0.7206491 0.7299145 0.7383258 0.7457353
[8] 0.7475100 0.7475856

```

Na matriz de valores lógicos, cada linha corresponde a uma cardinalidade (número de variáveis) do subconjunto preditor, e cada coluna corresponde a uma das variáveis predictoras. As colunas que tenham o valor lógico `TRUE`, na linha correspondente a  $k$  preditores, correspondem a variáveis que pertencem ao melhor subconjunto de  $k$  preditores. Repare-se como o melhor subconjunto de quatro preditores é o subconjunto `x1`, `x2`, `x6` e `x9`, escolhido pelo algoritmo de exclusão sequencial (nas suas duas versões). Aliás, em todos os passos intermédios do algoritmo, o subconjunto de  $k$  preditores escolhido acaba por revelar-se o subconjunto óptimo, ou seja, o subconjunto de preditores que está associado aos maiores valores do Coeficiente de Determinação.

(f) O ajustamento pedido nesta alínea produziu os seguintes resultados:

```

> summary(lm(I(y*0.06277) ~ x1 + I(x2*25.4) + I(x6*25.4) + I(5/9*(x9-32)), data=milho))
[...]
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) 3.5114712 1.5019053 2.338 0.0268 *
x1 0.0534744 0.0094015 5.688 4.25e-06 ***
I(x2 * 25.4) 0.0014800 0.0007877 1.879 0.0707 .

```

```

I(x6 * 25.4) 0.0058354 0.0028055 2.080 0.0468 *
I(5/9 * (x9 - 32)) -0.1102213 0.0599066 -1.840 0.0764 .

```

```

```

```
Residual standard error: 0.4677 on 28 degrees of freedom
```

```
Multiple R-squared: 0.7206, Adjusted R-squared: 0.6807
```

```
F-statistic: 18.06 on 4 and 28 DF, p-value: 1.954e-07
```

Comparando esta listagem com os resultados do modelo final produzido pelo algoritmo de exclusão sequencial, nas unidades de medida originais (ver alínea 19e), constata-se que as quantidades associadas à qualidade do ajustamento global ( $R^2$ , valor da estatística  $F$  no teste de ajustamento global) mantêm-se inalteradas. Trata-se dum consequência do facto de que as transformações de variáveis foram todas transformações lineares (afins). No entanto, e tal como sucedia na RLS, os valores das estimativas  $b_j$  são diferentes. O facto de que a informação relativa aos testes a  $\beta_j = 0$  se manter igual, para os coeficientes  $\beta_j$  que multiplicam as variáveis predictoras (isto é, quando  $j > 0$ ), sugere que se trata de alterações que apenas visam adaptar as estimativas às novas unidades de medida, não alterando globalmente o ajustamento.

20. (a) i. **Hipóteses:**  $H_0 : \beta_1 = \beta_2 = 0$ , vs.  $H_1 : \beta_1 \neq 0$  ou  $\beta_2 \neq 0$ .

**Estatística do teste:**  $F = \frac{n-(p+1)}{p} \frac{R^2}{1-R^2} \cap F_{(p,n-(p+1))}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica (Unilateral direita):** Rejeitar  $H_0$  se  $F_{calc} > f_{\alpha(p,n-(p+1))} = f_{0.05(2,28)} \approx 3.33$  (entre 3.32 e 3.39, nas tabelas).

**Conclusões:** O enunciado indica que o valor calculado da estatística é  $F_{calc} = 255$ .

Assim, *rejeita-se*  $H_0$ , indicando que o modelo RLM difere significativamente do modelo nulo.

- ii. Nos testes a que o coeficiente  $\beta_j$  de cada preditor ( $j = 1, 2$ ) seja nulo, os valores de prova dados no enunciado indicam que ambos são inferiores a  $\alpha = 0.05$ , pelo que haverá rejeição de  $H_0 : \beta_j = 0$  em ambos os casos e, ao nível  $\alpha = 0.05$ , qualquer das regressões lineares simples possíveis terá uma qualidade de ajustamento significativamente pior. Já ao nível  $\alpha = 0.01$  a situação é diferente. Enquanto o *p-value* para o teste a  $H_0 : \beta_1 = 0$  é  $p < 2 \times 10^{-16}$ , ou seja, indistinguível de zero e portanto indicando com grande convicção que  $\beta_1 \neq 0$ , já o valor de prova no teste a  $H_0 : \beta_2 = 0$  é  $p = 0.0145$  e portanto superior a  $\alpha = 0.01$ . Assim, e embora por pouco, não se rejeita a hipótese  $H_0 : \beta_2 = 0$  ao nível de significância  $\alpha = 0.01$ . Como tal, uma regressão linear simples de **Volume** sobre **Diametro** não difere significativamente (para  $\alpha = 0.01$ ) da regressão com dois preditores ajustada no enunciado.
- iii. Sabemos que numa regressão linear simples, o coeficiente de determinação é o quadrado do coeficiente de correlação entre o preditor e a variável resposta. Com base na matriz de correlações disponível no enunciado geral, temos que, na RLS de **Volume** sobre **Diametro** o coeficiente de determinação é  $R^2 = 0.9671194^2 = 0.9353199$ , enquanto que na RLS de **Volume** sobre **Altura** o coeficiente de determinação é  $R^2 = 0.5982497^2 = 0.3579027$ . Estes valores são coerentes com os resultados da alínea anterior. Quanto aos valores das estatísticas  $F$  nos testes de ajustamento global, podem ser obtidos pela fórmula da RLS,  $F = (n-2) \frac{R^2}{1-R^2}$ . Os valores nas duas regressões lineares simples são (e indicando o preditor pela sua inicial)  $F_D = 29 \times \frac{0.9353199}{1-0.9353199} = 419.3605$  e  $F_A = 29 \times \frac{0.3579027}{1-0.3579027} = 16.16449$ .

(b) Consideremos agora o modelo com base nas transformações logarítmicas das três variáveis originais. Designaremos por  $y$  o log-volume, por  $x_1$  o log-diâmetro e por  $x_2$  a log-altura.

i. Partindo da relação linear entre as variáveis logaritmizadas, tem-se:

$$\begin{aligned} \ln(y) = b_0 + b_1 \ln x_1 + b_2 \ln x_2 &\Leftrightarrow y = e^{b_0 + b_1 \ln x_1 + b_2 \ln x_2} \\ &\Leftrightarrow y = e^{b_0} e^{b_1 \ln x_1} e^{b_2 \ln x_2} \\ &\Leftrightarrow y = \underbrace{e^{b_0}}_{=b_0^*} e^{\ln x_1^{b_1}} e^{\ln x_2^{b_2}} \\ &\Leftrightarrow y = b_0^* x_1^{b_1} x_2^{b_2} . \end{aligned}$$

Assim,  $y$  é proporcional ao produto de potências de cada um dos preditores. A superfície em  $R^3$  ajustada à nuvem de pontos das observações originais terá, tendo em conta os valores disponíveis no enunciado, equação  $y = e^{-6.63162} x_1^{1.98265} x_2^{1.11712}$ , ou seja,  $Volume = 0.001318 \text{ Diâmetro}^{1.98265} \text{ Altura}^{1.11712}$ .

ii. Esta frase baseia-se numa comparação errada, uma vez que as escalas da variável resposta  $y$  (usadas para medir, resíduos e todas as Somas de Quadrados numa regressão, logo também usadas para obter os coeficientes de determinação e portanto também o valor da estatística  $F$ ) são diferentes nos dois modelos ajustados. Enquanto que na alínea anterior o volume era medido na escala original, nesta alínea a regressão linear usa a escala logarítmica para os volumes. Assim, o  $R^2$  da alínea anterior mede a proporção da variabilidade *dos volumes* observados que era explicada pela regressão então usada, nesta alínea o  $R^2$  mede a variabilidade *dos log-volumes* observados que é explicada pela nova regressão. Os  $SQT$ s de cada alínea não são iguais. Não são correctas as comparações referidas na frase do enunciado.

(c) A troca de variável resposta piorou claramente o valor de  $R^2$  do ajustamento. Este resultado pode parecer surpreendente à primeira vista, uma vez que do ponto de vista algébrico, uma relação da forma  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$  é equivalente a  $x_2 = \frac{y - \beta_0 - \beta_1 x_1}{\beta_2} = \beta_0^* + \beta_1^* x_1 + \beta_2^* y$  (com  $\beta_0^* = \frac{-\beta_0}{\beta_2}$ ,  $\beta_1^* = \frac{-\beta_1}{\beta_2}$  e  $\beta_2^* = \frac{1}{\beta_2}$ ). Além disso, numa regressão linear simples, a troca do preditor e da variável resposta, se bem que muda a equação da recta ajustada, não muda a qualidade do ajustamento (uma vez que  $R^2 = r_{xy}^2$ , e o coeficiente de correlação é simétrico nos seus argumentos). Mas numa regressão linear múltipla, permutar a variável resposta com um dos preditores pode, como este exemplo ilustra, gerar um modelo de qualidade bastante diferente. O exemplo sugere a razão de ser deste facto: as variáveis **Volume** e **Diâmetro** estão fortemente correlacionadas entre si. Qualquer modelo de regressão linear que tenha uma dessas variáveis como variável resposta, e a outra como preditor, terá de ter  $R^2 \geq (0.9671194)^2 = 0.9353199$ . Mas a variável **Altura**, que foi agora colocada como variável resposta, não está fortemente correlacionada com nenhuma das duas outras. Ao desempenhar o papel de variável resposta, com as outras duas variáveis como preditores, o valor do  $R^2$  resultante poderá ser elevado, mas como este exemplo ilustra, poderá não o ser.

21. Vamos contruir o intervalo de confiança a  $(1 - \alpha) \times 100\%$  para  $\mathbf{a}^t \vec{\beta}$ , a partir da distribuição indicada no enunciado. Sendo  $t_{\frac{\alpha}{2}}$  o valor que, numa distribuição  $t_{n-(p+1)}$ , deixa à direita uma região de probabilidade  $\alpha/2$ , temos a seguinte afirmação probabilística, na qual trabalhamos a dupla desigualdade até deixar a combinação linear (para a qual se quer o intervalo de confiança)

isolada no meio:

$$\begin{aligned}
 P \left[ -t_{\frac{\alpha}{2}} < \frac{\mathbf{a}^t \vec{\hat{\beta}} - \mathbf{a}^t \vec{\beta}}{\hat{\sigma}_{\mathbf{a}^t \vec{\beta}}} < t_{\frac{\alpha}{2}} \right] &= 1 - \alpha \\
 P \left[ -t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\mathbf{a}^t \vec{\beta}} < \mathbf{a}^t \vec{\hat{\beta}} - \mathbf{a}^t \vec{\beta} < t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\mathbf{a}^t \vec{\beta}} \right] &= 1 - \alpha \\
 P \left[ t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\mathbf{a}^t \vec{\beta}} > \mathbf{a}^t \vec{\beta} - \mathbf{a}^t \vec{\hat{\beta}} > -t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\mathbf{a}^t \vec{\beta}} \right] &= 1 - \alpha \quad (\text{multiplicando por } -1) \\
 P \left[ \mathbf{a}^t \vec{\hat{\beta}} - t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\mathbf{a}^t \vec{\beta}} < \mathbf{a}^t \vec{\beta} < \mathbf{a}^t \vec{\hat{\beta}} + t_{\frac{\alpha}{2}} \cdot \hat{\sigma}_{\mathbf{a}^t \vec{\beta}} \right] &= 1 - \alpha
 \end{aligned}$$

Assim, calculando o valor  $\mathbf{a}^t \mathbf{b} = a_0 b_0 + a_1 b_1 + \dots + a_p b_p$  do estimador  $\mathbf{a}^t \vec{\hat{\beta}}$  e o erro padrão  $\hat{\sigma}_{\mathbf{a}^t \vec{\beta}}$ , para a nossa amostra, temos o intervalo a  $(1-\alpha) \times 100\%$  de confiança para  $\mathbf{a}^t \vec{\beta} = a_0 \beta_0 + a_1 \beta_1 + \dots + a_p \beta_p$ :

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

22. Parte-se duma regressão linear simples relacionando a variável resposta **Peso** e o preditor **Calibre**.

- (a) A ordenada na origem natural é  $\beta_0 = 0$ : a calibre nulo corresponde inexistência de fruto, ou seja, peso nulo. O intervalo a 95% de confiança para a ordenada na origem é dado por:

$$\left] b_0 - t_{\frac{\alpha}{2}(n-2)} \cdot \hat{\sigma}_{\hat{\beta}_0} \quad , \quad b_0 + t_{\frac{\alpha}{2}(n-2)} \cdot \hat{\sigma}_{\hat{\beta}_0} \quad [$$

No enunciado verifica-se que  $b_0 = -210.3137$ , com erro padrão associado  $\hat{\sigma}_{\hat{\beta}_0} = 3.8078$ . Tem-se ainda  $t_{0.025(1271)} \approx 1.96$ . Logo, o IC pedido é  $] -217.777, -202.8504 [$ . Este intervalo está muito longe de incluir o valor natural  $\beta_0 = 0$ , pelo que essa eventualidade pode ser excluída. Não sendo um resultado encorajador, a verdade é que não faz sentido utilizar um modelo deste tipo para frutos de calibre próximo de zero. Os calibres utilizados no ajustamento do modelo variaram entre 53 e 79, pelo que deve evitar-se utilizar este modelo para calibres muito distantes da gama de calibres observados.

- (b) Nesta alínea ajustou-se um polinómio de segundo grau, através dum modelo de regressão múltipla em que  $X_1 = \text{Calibre}$  e  $X_2 = \text{Calibre}^2$ . A equação de base neste modelo é  $\text{Peso} = \beta_0 + \beta_1 \text{Calibre} + \beta_2 \text{Calibre}^2$ .

- i. A equação da parábola ajustada é:  $\text{Peso} = 72.33140 - 3.38747 \text{Calibre} + 0.06469 \text{Calibre}^2$ . Observe como a ordenada na origem e o coeficiente da variável **Calibre** são radicalmente diferentes do que eram na regressão linear simples.
- ii. O modelo linear e o modelo quadrático são equivalentes caso  $\beta_2 = 0$ . Essa hipótese pode ser testada como qualquer outro teste  $t$  a um parâmetro  $\beta_j$  individual do modelo:

**Hipóteses:**  $H_0 : \beta_2 = 0$  vs.  $H_1 : \beta_2 \neq 0$ .

**Estatística do teste:**  $T = \frac{\hat{\beta}_2 - 0}{\hat{\sigma}_{\hat{\beta}_2}} \cap t_{n-(p+1)}$

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica (Bilateral):** Rejeitar  $H_0$  se  $|T_{\text{calc}}| > t_{\alpha/2(n-3)} = t_{0.025(1270)} \approx 1.962$ .



**Conclusões:** O valor calculado da estatística do teste é dado no enunciado, na penúltima coluna da tabela *Coefficients*:  $T_{calc} = \frac{0.06469}{0.01067} = 6.064$ . Logo, rejeita-se claramente a hipótese nula  $\beta_2 = 0$ , pelo que o modelo polinomial (quadrático) tem um ajustamento significativamente melhor que o modelo linear. Repare-se como este resultado está associado a um aumento bastante pequeno do coeficiente de determinação  $R^2$  (de 0.8638 para 0.8677). Este facto está, mais uma vez, associado à grande dimensão da amostra ( $n = 1273$ ), que permite considerar significativas diferenças tão pequenas quanto estas.

23. (a) O modelo de regressão linear múltipla relaciona uma variável resposta  $Y$  com  $p$  variáveis preditoras  $X_1, X_2, \dots, X_p$ . Designando por  $\vec{Y}$  o vector das  $n$  observações da variável resposta  $Y$ ,  $\vec{\epsilon}$  o vector dos  $n$  erros aleatórios correspondentes,  $\vec{\beta}$  o vector dos  $p + 1$  parâmetros do modelo,  $\beta_0, \beta_1, \dots, \beta_p$ , e  $\mathbf{X}$  a matriz  $n \times (p + 1)$ , cuja primeira coluna é constituída por  $n$  uns e cada uma das restantes  $p$  colunas contém as  $n$  observações duma variável preditora, tem-se:

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

O modelo de regressão linear múltipla é então dado por:

- i.  $\vec{Y} = \mathbf{X}\vec{\beta} + \vec{\epsilon}$
- ii.  $\vec{\epsilon} \cap \mathcal{N}_n(\vec{0}, \sigma^2\mathbf{I}_n)$ ,

sendo  $\vec{0}$  o vector de  $n$  zeros e  $\mathbf{I}_n$  a matriz identidade  $n \times n$ . Na segunda condição, indica-se que o vector dos erros aleatórios segue uma distribuição Multinormal, com vector médio dado pelo vector de zeros (ou seja, cada erro aleatório individual tem valor esperado zero) e matriz de variâncias-covariâncias diagonal, com os elementos diagonais todos iguais a  $\sigma^2$ . Uma vez que, numa matriz de (co-)variâncias os elementos diagonais representam as variâncias de cada componente do vector, esta condição indica que  $V[\epsilon_i] = \sigma^2, \forall i$ . O facto de os elementos não diagonais da matriz  $\sigma^2\mathbf{I}_n$  serem todos nulos equivale a dizer que a covariância entre elementos diferentes do vector aleatório dos erros é sempre nula (ou seja,  $Cov[\epsilon_i, \epsilon_j] = 0$ , sempre que  $i \neq j$ ) e, como sabemos, numa distribuição Multinormal tal facto implica a independência desses elementos.

- (b) O vector  $\vec{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)^t$  dos estimadores dos  $p + 1$  parâmetros dum modelo linear é dado (ver formulário) por  $\vec{\beta} = (\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t\vec{Y}$ . Mas, pelo modelo, tem-se  $\vec{Y} = \mathbf{X}\vec{\beta} + \vec{\epsilon}$ . Substituindo, tem-se:

$$\vec{\beta} = (\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t(\mathbf{X}\vec{\beta} + \vec{\epsilon}) = \underbrace{(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t\mathbf{X}}_{=\mathbf{I}}\vec{\beta} + (\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t\vec{\epsilon} = \vec{\beta} + (\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t\vec{\epsilon},$$

como se pedia para mostrar.

- (c) A expressão da alínea anterior é a soma dum vector não aleatório,  $\vec{\beta}$ , com um vector aleatório,  $(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t\vec{\epsilon}$ . Ora, para qualquer vector aleatório  $\vec{W}$  e vector não aleatório  $\mathbf{a}$  verifica-se  $E[\vec{W} + \mathbf{a}] = E[\vec{W}] + \mathbf{a}$ . Logo, no nosso caso, tem-se:  $E[\vec{\beta}] = E[\vec{\beta} + (\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t\vec{\epsilon}] =$

$\vec{\beta} + E[(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \vec{\epsilon}]$ . A segunda parcela é o vector esperado dum vector que resulta de multiplicar uma matriz não aleatória ( $(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t$ ) por um vector aleatório ( $\vec{\epsilon}$ ). Por outra propriedade operatória dos vectores esperados, tem-se  $E[\mathbf{B} \vec{\mathbf{W}}] = \mathbf{B} E[\vec{\mathbf{W}}]$ , onde  $\mathbf{B}$  é uma matriz não aleatória. Assim,  $E[\vec{\beta}] = \vec{\beta} + E[(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \vec{\epsilon}] = \vec{\beta} + (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \underbrace{E[\vec{\epsilon}]}_{=\vec{0}} = \vec{\beta} + \vec{0} = \vec{\beta}$ .

Por outro lado, tendo em conta a propriedade operatória geral de matrizes de (co-)variâncias,  $V[\vec{\mathbf{W}} + \mathbf{a}] = V[\vec{\mathbf{W}}]$ , tem-se  $V[\vec{\beta}] = V[\vec{\beta} + (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \vec{\epsilon}] = V[(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \vec{\epsilon}]$ . Outra propriedade operatória de matrizes de (co-)variâncias diz-nos que  $V[\mathbf{B} \vec{\mathbf{W}}] = \mathbf{B} V[\vec{\mathbf{W}}] \mathbf{B}^t$ , para uma matriz não aleatória  $\mathbf{B}$ . Logo (e sendo no nosso caso  $\mathbf{B} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t$ ), tem-se:

$$\begin{aligned} V[\vec{\beta}] &= V[(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \vec{\epsilon}] = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t V[\vec{\epsilon}] [(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t]^t \\ &= (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \sigma^2 \mathbf{I}_n \mathbf{X} [(\mathbf{X}^t \mathbf{X})^{-1}]^t = \sigma^2 (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{X} [(\mathbf{X}^t \mathbf{X})^{-1}]^t = \sigma^2 (\mathbf{X}^t \mathbf{X})^{-1} . \end{aligned}$$

24. O  $R^2$  modificado definiu-se como  $R_{mod}^2 = 1 - \frac{QMRE}{QMT}$  onde  $QMT = \frac{SQT}{n-1} = s_y^2$ .

(a) Tem-se  $R_{mod}^2 = 1 - \frac{QMRE}{QMT} = 1 - \frac{\frac{SQRE}{n-(p+1)}}{\frac{SQT}{n-1}} = 1 - \frac{n-1}{n-(p+1)} \times \frac{SQRE}{SQT}$ . Mas  $\frac{SQRE}{SQT} = \frac{SQT - SQR}{SQT} = 1 - R^2$ .

Substituindo na expressão anterior, tem-se o resultado pretendido:  $R_{mod}^2 = 1 - (1 - R^2) \frac{n-1}{n-(p+1)}$ .

(b) Por definição, a estatística do teste  $F$  de ajustamento global tem valor  $F_{calc} = \frac{QMR}{QMRE} = \frac{n-(p+1)}{p} \cdot \frac{R^2}{1-R^2}$ . Ora, usando a expressão da alínea anterior, tem-se  $R^2 - R_{mod}^2 = R^2 - 1 + (1 - R^2) \frac{n-1}{n-(p+1)} = (1 - R^2) \left[ -1 + \frac{n-1}{n-(p+1)} \right] = (1 - R^2) \left[ \frac{n-(p+1) + p - 1}{n-(p+1)} \right] = (1 - R^2) \frac{p}{n-(p+1)}$ . Logo,  $\frac{R^2}{R^2 - R_{mod}^2} = \frac{R^2}{(1 - R^2) \frac{p}{n-(p+1)}} = \frac{n-(p+1)}{p} \times \frac{R^2}{1 - R^2} = F_{calc}$ , como se queria mostrar.

(c) Usando os resultados da primeira alínea, tem-se  $R_{mod}^2 < 0 \Leftrightarrow 1 - (1 - R^2) \frac{n-1}{n-(p+1)} < 0 \Leftrightarrow 1 < (1 - R^2) \frac{n-1}{n-(p+1)} \Leftrightarrow \frac{n-(p+1)}{n-1} < 1 - R^2 \Leftrightarrow R^2 < 1 - \frac{n-(p+1)}{n-1} = \frac{n-1 - n + p + 1}{n-1} = \frac{p}{n-1}$ , como se pedia para mostrar. Se esta condição se verifica, tem-se, a partir da expressão da alínea anterior, que  $F_{calc}$  terá valor inferior a 1. Uma rápida vista de olhos pelas tabelas da distribuição  $F$  mostra que valores de  $F_{calc}$  inferiores a 1 nunca conduzem (para os níveis de significância usuais) à rejeição da  $H_0$ , pelo que o modelo ajustado não passa o teste de ajustamento global. Assim, ajustamentos de modelos em que  $p$  seja pouco menor do que  $n-1$  podem não passar o teste de ajustamento global, mesmo com valores relativamente elevados de  $R^2$ .

## 2 Análise de Variância

1. (a) Trata-se dum delineamento a um único factor (as variedades de tomate), sendo a variável resposta  $Y$  a resistência da película (em *gf*). Em cada um dos  $k=6$  níveis do factor há  $n_c=3$  repetições (as parcelas). O número igual de repetições nas 6 situações experimentais significa que o delineamento é equilibrado. O modelo ANOVA a um factor corresponde a:
  - i. A resistência  $Y_{ij}$ , na  $j$ -ésima parcela ( $j=1, 2, 3$ ) associada à variedade  $i$  ( $i=1, \dots, 6$ ), é dada por:

$$Y_{ij} = \mu_1 + \alpha_i + \epsilon_{ij} , \quad \forall i, j ,$$

sendo  $\mu_1$  a resistência esperada da primeira variedade;  $\alpha_i = \mu_i - \mu_1$  o efeito (acréscimo à resistência média da primeira variedade) da variedade  $i$  (com  $\alpha_1 = 0$ ); e  $\epsilon_{ij}$  o erro aleatório da observação  $Y_{ij}$ . Iremos (tal como o programa R) admitir que as variedades estão ordenadas por ordem alfabética, com os nomes de nível numéricos à cabeça, pelo que a primeira variedade acima referida é a variedade 18.

- ii. Admite-se que os erros aleatórios são todos Normais, de média zero e variâncias homogêneas, ou seja, para qualquer  $i, j$ :

$$\epsilon_{ij} \cap \mathcal{N}(0, \sigma^2) .$$

- iii. Admite-se que os erros aleatórios  $\epsilon_{ij}$  são independentes.

- (b) A tabela-resumo terá apenas duas linhas (além da linha correspondente aos Totais), associadas respectivamente aos efeitos do Factor e à variabilidade Residual.

- i. Sabemos que os graus de liberdade dos efeitos do factor são  $k-1=5$  e que os graus de liberdade residuais são  $n-k=18-6=12$ . As fórmulas para as Somas de Quadrados são dadas no formulário. A Soma de Quadrados Residual é  $SQRE = \sum_{i=1}^k (n_i - 1)s_i^2$  e, tratando-se dum delineamento equilibrado com  $n_c = 3$  repetições em cada nível, tem-se  $SQRE = (n_c - 1) \sum_{i=1}^k s_i^2$ . Usando as variâncias amostrais de nível dadas no enunciado, vem  $SQRE = 2 \times (14713.08 + 367.9434 + 5881.921 + 33132.64 + 5.414433 + 47.11163) = 108\,296.2$ . É possível calcular  $SQF$  através da sua fórmula, uma vez que são disponibilizadas as médias amostrais de nível e globais. Mas é mais simples obter esse valor constatando que, numa ANOVA a um factor, se tem  $SQF = SQT - SQRE$ . No nosso caso  $SQT = (n-1)s_y^2 = 17 \times 34\,517.82 = 586\,802.9$ . Logo,  $SQF = 478\,506.7$ . Dividindo estas Somas de Quadrados pelos graus de liberdade antes referidos obtêm-se os Quadrados Médios, e dividindo  $QMF$  por  $QMRE$  obtém-se o valor calculado da estatística do teste  $F$  aos efeitos do factor. Eis a tabela-resumo:

|          | g.l. | SQs       | Quadrados Médios                     | $F_{calc}$                                                              |
|----------|------|-----------|--------------------------------------|-------------------------------------------------------------------------|
| Factor   | 5    | 478 506.7 | $\frac{478\,506.7}{5} = 95\,701.35$  | $F_{calc} = \frac{QMF}{QMRE} = \frac{95\,701.35}{9\,024.685} = 10.6044$ |
| Residual | 12   | 108 296.2 | $\frac{108\,296.2}{12} = 9\,024.685$ |                                                                         |

- ii. Usando o R, confirmamos a tabela-resumo agora obtida:

```
> tomate.aov <- aov(res.pel ~ variedade , data=tomate)
> summary(tomate.aov)
 Df Sum Sq Mean Sq F value Pr(>F)
variedade 5 478507 95701 10.6 0.000448
Residuals 12 108296 9025
```

- (c) Eis o teste aos efeitos do factor (variedade):

**Hipóteses:**  $H_0 : \alpha_i = 0, \forall i$  vs.  $H_1 : \exists i$  tal que  $\alpha_i \neq 0$ .

**Estatística do Teste:**  $F = \frac{QMF}{QMRE} \cap F_{[k-1, n-k]}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica:** (Unilateral direita) Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(5,12)} = 3.11$ .

**Conclusões:** Como  $F_{calc} = 10.6044 > 3.11$ , rejeita-se  $H_0$ , concluindo-se que existem de efeitos de variedade (ao nível  $\alpha = 0.05$ ), o que corresponde a afirmar que existem variedades de tomate cujas películas têm resistência média diferentes de outras.

- (d) O valor de prova (*p-value*) associado ao valor calculado da estatística de teste é  $p=0.000448$ . Pela própria definição de *p-value*, esta é a área à direita de  $F_{calc}=10.6044$ , numa distribuição  $F_{[5,12]}$ . Logo, seria preciso fazer um teste de hipóteses com nível de significância  $\alpha=0.000448$  (ou inferior) para que  $F_{calc}$  não pertencesse à Região Crítica e a conclusão do teste pudesse ser a de não rejeitar  $H_0$ .
- (e) Tal como nas regressões lineares, a primeira coluna da matriz  $\mathbf{X}$  é uma coluna de uns. No contexto duma ANOVA a um factor, as restantes colunas são variáveis indicatrizes de pertença de cada observação a um dos níveis do factor, ou seja, colunas com apenas dois valores: “1” associado a observações que pertencem ao nível do factor em causa, e “0” associado a observações associadas a outros níveis do factor. A restrição imposta no modelo ( $\alpha_1=0$ ) implica que não há indicatriz do primeiro nível do factor, neste caso, o nível “18”. Assim, neste caso teremos uma primeira coluna de  $n=18$  uns e cinco colunas indicatrizes dos segundo, terceiro, quarto, quinto e sexto níveis do factor ( $\mathcal{I}_2, \mathcal{I}_3, \mathcal{I}_4, \mathcal{I}_5$  e  $\mathcal{I}_6$ ), como se pode confirmar através do comando referido no enunciado:

```
> model.matrix(tomate.aov)
 (Intercept) variedade28 variedade29 variedade40C variedadeAce variedadeRoma
1 1 0 0 0 0 0
2 1 0 0 0 0 0
3 1 0 0 0 0 0
4 1 1 0 0 0 0
5 1 1 0 0 0 0
6 1 1 0 0 0 0
7 1 0 1 0 0 0
8 1 0 1 0 0 0
9 1 0 1 0 0 0
10 1 0 0 1 0 0
11 1 0 0 1 0 0
12 1 0 0 1 0 0
13 1 0 0 0 0 1
14 1 0 0 0 0 1
15 1 0 0 0 0 1
16 1 0 0 0 1 0
17 1 0 0 0 1 0
18 1 0 0 0 1 0
```

A ordem dos níveis do factor no R é, por omissão, a ordem alfabética dos nomes dos níveis. Mas essa pode não ser a ordem pela qual as observações surgem nas linhas da *data frame* com os dados. Neste exemplo, a variedade **Roma** surge como último nível (última coluna de  $\mathbf{X}$ ), mas as observações dessa variedade não estão nas linhas finais da *data frame*, razão pela qual as duas colunas finais de  $\mathbf{X}$  parecem 'trocadas'.

- (f) Os valores ajustados  $\hat{Y}_{ij}$ , numa ANOVA a um factor, são as médias amostrais do nível a que cada observação pertence. Assim, tem-se:

```
> fitted(tomate.aov)
 1 2 3 4 5 6 7 8
560.6433 560.6433 560.6433 241.4833 241.4833 241.4833 290.9500 290.9500
 9 10 11 12 13 14 15 16
290.9500 705.7800 705.7800 705.7800 332.1067 332.1067 332.1067 377.2533
 17 18
377.2533 377.2533
```

Estas são as médias de variedade dadas no enunciado (arredondadas a uma casa decimal).

(g) O facto dos resíduos se encontrarem ‘empilhados’ em seis colunas é o reflexo natural do facto, referido na alínea anterior, de apenas haver seis diferentes valores ajustados nesta ANOVA: as seis médias amostrais de cada variedade,  $\hat{y}_{ij} = \bar{y}_i$  ( $j = 1, 2, 3$ ). Este facto ajuda a identificar as observações associadas aos resíduos de maior magnitude. Assim, por exemplo, o maior resíduo (em módulo) corresponde ao ponto no canto inferior direito. Por estar associado a uma média  $\bar{y}_i$  de aproximadamente 700, tem de corresponder à variedade 40C. Por ser um resíduo negativo, tem de corresponder a uma observação com valor inferior à média dessa variedade, o que apenas acontece com a primeira das três observações desse nível. Assim, a observação a que corresponde o referido resíduo é a observação  $y_{4,1} = 503.51$ . Embora o número de repetições em cada nível ( $n_c = 3$ ) seja muito baixo, e portanto susceptível de gerar impressões enganadoras, o gráfico sugere alguma heterogeneidade nas variâncias de  $Y_{ij}$  em cada nível. Os valores das variâncias amostrais de nível indicam que há variedades com muito pouca variabilidade nas resistências observadas (como a *Ace*, com  $s_5^2 = 5.414433$ ) e outras com uma variabilidade muito maior (como a *29*, com  $s_3^2 = 5881.921$ , mais de cem vezes maior).

2. Neste exercício sobre os estomas das folhas de café, não estão disponíveis os dados originais. Apenas se conhece a tabela dos valores médios e variâncias amostrais de cada variedade.

(a) A variável resposta  $Y$  é o comprimento médio dos estomas das folhas duma planta. Para explicar a variabilidade dos valores desta variável, apenas se dispõe de um factor: o factor variedade, com  $k = 3$  níveis (as três variedades indicadas no enunciado). O modelo ANOVA é assim o modelo a um factor, semelhante ao do primeiro exercício. É um delineamento equilibrado, pois existem  $n_i = 12$  observações para qualquer variedade ( $i = 1, 2, 3$ ), perfazendo um total de  $n = 3 \times 12 = 36$  observações  $Y_{ij}$ . Eis o modelo:

i.  $Y_{ij} = \mu_1 + \alpha_i + \epsilon_{ij}$ ,  $\forall i = 1, 2, 3$ ,  $j = 1, 2, \dots, 12$ , com  $\alpha_1 = 0$ , onde

- $Y_{ij}$  indica o comprimento médio dos estomas das folhas da planta  $j$  da variedade  $i$ ;
- $\mu_1$  indica o comprimento médio populacional dos estomas das folhas de plantas da primeira variedade ( $i = 1$ ) que é, por ordem alfabética, a variedade CA;
- $\alpha_i$  indica o efeito (acréscimo em relação à média da variedade CA) da variedade  $i$ ; e
- $\epsilon_{ij}$  indica o erro aleatório associado à observação  $Y_{ij}$ .

ii.  $\epsilon_{ij} \cap \mathcal{N}(0, \sigma^2)$ ,  $\forall i, j$ .

iii.  $\{\epsilon_{ij}\}_{i,j}$  constitui um conjunto de variáveis aleatórias independentes.

(b) Começemos pelo cálculo das Somas de Quadrados. Uma vez que o delineamento é equilibrado (igual número de observações em cada nível), a média global da totalidade das 36 observações ( $\bar{y}_{..}$ ) é a média simples das três médias de nível dadas na tabela:  $\bar{y}_{..} = (22.85833 + 19.49333 + 25.31583)/3 = 22.55583$ . Tendo em conta as fórmulas vistas nas aulas teóricas e os valores dados no enunciado, temos:

$$SQRE = (n_c - 1) \sum_{i=1}^3 s_i^2 = 11 \times (13.69303 + 2.725424 + 9.388936) = 284.1983 ;$$

$$\begin{aligned} SQF &= n_c \sum_{i=1}^3 (\bar{y}_i - \bar{y}_{..})^2 \\ &= 12 \times ((22.85833 - 22.55583)^2 + (19.49333 - 22.55583)^2 + (25.31583 - 22.55583)^2) \\ &= 205.0561, \end{aligned}$$

Logo, tem-se a seguinte tabela-resumo:

| Fonte    | g.l.       | SQ                | QM                                  | $F_{calc}$                    |
|----------|------------|-------------------|-------------------------------------|-------------------------------|
| Factor   | $k-1 = 2$  | $SQF = 205.0561$  | $QMF = \frac{SQF}{k-1} = 102.5281$  | $\frac{QMF}{QMRE} = 11.90516$ |
| Resíduos | $n-k = 33$ | $SQRE = 284.1983$ | $QMRE = \frac{SQRE}{n-k} = 8.61207$ |                               |

- (c) Neste caso, e uma vez que não são conhecidas as observações individuais, apenas é possível calcular a variância da totalidade das  $n = 36$  observações recorrendo à decomposição da Soma de Quadrados Total correspondente a esta ANOVA:

$$s_y^2 = \frac{SQT}{n-1} = \frac{SQF + SQRE}{n-1} = \frac{205.0561 + 284.1983}{35} = \frac{489.2544}{35} = 13.9787 .$$

Repare-se que este valor *não é* a média das variâncias amostrais de nível.

- (d) Embora se possa escrever as hipóteses do teste com base nos efeitos  $\alpha_i$  do factor (como se fez no exercício anterior), nas ANOVAs a um único factor é equivalente formular as hipóteses em termos das médias populacionais (valores esperados das observações  $E[Y_{ij}] = \mu_i = \mu_1 + \alpha_i$ ) em cada nível do factor. Eis o teste com  $\alpha = 0.05$ :

**Hipóteses:**  $H_0 : \mu_1 = \mu_2 = \mu_3$  vs.  $H_1 : \exists i, i'$  tal que  $\mu_i \neq \mu_{i'}$ .

**Estatística do teste:**  $F = \frac{QMF}{QMRE} \cap F_{(k-1, n-k)}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica (Unilateral Direita):** Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(2,33)} \approx 3.30$  (entre os valores tabelados 3.23 e 3.32).

**Conclusões:** O valor da estatística do teste foi calculado na alínea anterior:  $F_{calc} = 11.90516$ . É um valor significativo ao nível  $\alpha = 0.05$  e rejeita-se  $H_0$  a favor da hipótese de que existem efeitos do factor, ou seja, de que o comprimento médio dos estomas das folhas não é igual em todas as variedades.

O valor de prova associado à estatística calculada é (tendo em conta a natureza unilateral direita do teste)  $P[F_{(2,33)} > F_{calc}] = P[F_{(2,33)} > 11.90516]$ . Não é possível obter este valor nas tabelas, mas pode calcular-se essa probabilidade com o auxílio do **R**:

$> 1-pf(11.90516, 2, 33)$

[1] 0.000128065

Assim, tem-se  $p = 0.000128065$ .

- (e) **[Material Complementar]** Sabemos que duas médias de nível  $\mu_i$  e  $\mu_{i'}$  devem ser consideradas diferentes caso as respectivas médias amostrais difiram (em módulo) mais do que o termo de comparação  $q_{\alpha(k, n-k)} \sqrt{\frac{QMRE}{n_c}}$ , onde  $q_{\alpha(k, n-k)}$  corresponde ao valor que deixa à sua direita uma região de probabilidade  $\alpha$  numa distribuição de Tukey de parâmetros  $k$  e  $n-k$ , e  $n_c$  indica o número comum de observações em cada nível do factor (o resultado que sustenta o teste de Tukey parte do pressuposto que o delineamento é equilibrado). No nosso caso tem-se  $k = 3$  e  $n = 36$ . Trabalhando (como pedido no enunciado) com  $\alpha = 0.05$ , e recorrendo às tabelas da distribuição de Tukey (tabelas específicas, disponíveis na página *web* da disciplina), tem-se  $q_{0.05(3,33)} = 3.47$ . Um valor mais preciso pode ser obtido através do comando `qtukey` do **R**:

```
> qtukey(0.95, 3, 33)
[1] 3.470189
```

Sabemos pela alínea (b) que  $QMRE = 8.61207$  e também que  $n_c = 12$ . Logo, o termo de comparação é dado por  $q_{\alpha(k,n-k)} \sqrt{\frac{QMRE}{n_c}} = 3.470189 \times \sqrt{\frac{8.61207}{12}} = 2.490459$ . Calculando as diferenças entre as médias amostrais de cada variedade, obtém-se a seguinte tabela:

| $ \bar{y}_i - \bar{y}_{i'} $ | CA ( $i'=1$ ) | CL ( $i'=2$ ) | PR ( $i'=3$ ) |
|------------------------------|---------------|---------------|---------------|
| CA ( $i=1$ )                 | –             | 3.3650        | 2.4575        |
| CL ( $i=2$ )                 | 3.3650        | –             | 5.8225        |
| PR ( $i=3$ )                 | 2.4575        | 5.8225        | –             |

Assim, ao nível de significância  $\alpha = 0.05$ , o comprimento médio dos estomas de folhas da variedade CL é diferente, quer do comprimento médio da variedade CA, quer do comprimento médio da variedade PR. No entanto, não se pode considerar (por pouco) significativamente diferentes os comprimentos médios dos estomas de folhas das variedades CA e PR.

Existem duas formas frequentes de representar esta conclusão, sendo usual em ambas ordenar os níveis do factor por ordem crescente das respectivas médias, e:

- i. sublinhando-se com traços os grupos de níveis cujas médias não diferem significativamente o que, nesta alínea (ao nível  $\alpha=0.05$ ) produz o seguinte resultado:

| CL       | CA              | PR              |
|----------|-----------------|-----------------|
| 19.49333 | <u>22.85833</u> | <u>25.31583</u> |

- ii. ou colocando uma mesma letra ao lado das variedades cujas médias não se consideram significativamente diferentes, por exemplo:

| CL                    | CA                    | PR                    |
|-----------------------|-----------------------|-----------------------|
| 19.49333 <sup>a</sup> | 22.85833 <sup>b</sup> | 25.31583 <sup>b</sup> |

Assim, a média de CL é significativamente diferente das médias, quer de CA, quer de PR (com quem não partilha letras em comum), mas já a média da variedade CA não difere significativamente da média de PR (uma vez que partilham a mesma letra).

3. A variável resposta  $Y$  é, neste caso, a variação de massa (coluna `variacao.massa` na `data frame`). Existem ao todo  $n = 50$  observações.

- (a) Para estudar este problema através duma ANOVA, ignora-se os valores numéricos das concentrações de dióxido de carbono, tratando cada diferente concentração apenas como um diferente tratamento. Assim, o factor  $CO_2$  terá  $k=5$  níveis, havendo ( $n_i=10=n_c$ ) observações para cada concentração de  $CO_2$  (nível do factor). O modelo ANOVA associado a este delineamento é o seguinte:

- i.  $Y_{ij} = \mu_1 + \alpha_i + \epsilon_{ij}$ ,  $\forall i = 1, 2, 3, 4, 5$ ,  $j = 1, 2, \dots, 10$ , com  $\alpha_1 = 0$ , onde
  - $Y_{ij}$  indica a variação de massa para a  $j$ -ésima repetição associada à  $i$ -ésima concentração de  $CO_2$ ;
  - $\mu_1$  indica o variação de massa média (populacional) na ausência de  $CO_2$  ( $i = 1$ );
  - $\alpha_i$  indica o efeito (acréscimo em relação à média populacional do primeiro nível) da  $i$ -ésima concentração de dióxido de carbono, isto é,  $\alpha_i = \mu_i - \mu_1$ ; e
  - $\epsilon_{ij}$  indica o erro aleatório associado à observação  $Y_{ij}$ .
- ii.  $\epsilon_{ij} \cap \mathcal{N}(0, \sigma^2)$ ,  $\forall i, j$ .
- iii.  $\{\epsilon_{ij}\}_{i,j}$  constitui um conjunto de variáveis aleatórias independentes.

Constrói-se a tabela-resumo da ANOVA com o auxílio do R (dados disponíveis na *data frame* C02, com os valores da variável resposta na coluna `variacao.massa` e os níveis de  $CO_2$  no factor `C02.factor`):

```
> summary(aov(variacao.massa ~ C02.factor, data=C02))
 Df Sum Sq Mean Sq F value Pr(>F)
C02.factor 4 11274 2818.6 101.6 <2e-16 ***
Residuals 45 1248 27.7
```

O teste  $F$  desta ANOVA diz respeito à possível existência de efeitos do Factor, ou seja,

**Hipóteses:**  $H_0 : \alpha_i = 0, \forall i = 2, 3, 4, 5$  vs.  $H_1 : \exists i = 2, 3, 4, 5$  tal que  $\alpha_i \neq 0$ .

**Estatística do teste:**  $F = \frac{QMF}{QMRE} \sim F_{(k-1, n-k)}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica (Unilateral Direita):** Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(4,45)} \approx 2.58$ .

**Conclusões:** O valor da estatística do teste é  $F_{calc} = 101.6$ , um valor claramente significativo ao nível  $\alpha = 0.05$ . Rejeita-se  $H_0$  a favor da hipótese de que existem efeitos do Factor, ou seja, que as concentrações de  $CO_2$  estão associadas a diferentes variações médias na massa das culturas do *Pseudomonas fragi*.

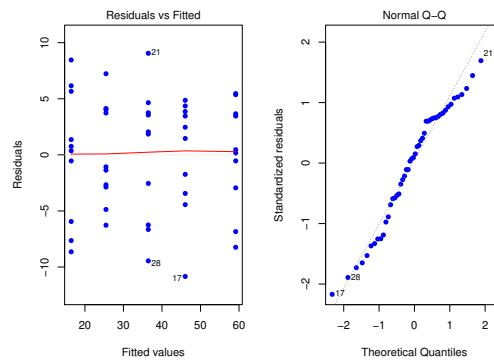
Como em qualquer modelo linear, o resíduo é a diferença entre cada valor observado da variável resposta e o correspondente valor ajustado pelo modelo, ou seja, e usando a notação da ANOVA a 1 Factor,  $e_{ij} = y_{ij} - \hat{y}_{ij}$ . Sabe-se que, num modelo ANOVA a um factor, o valor ajustado numa dada observação corresponde à média amostral das observações no mesmo nível do factor:  $\hat{y}_{ij} = \bar{y}_i$ . Assim, todas as observações do primeiro grupo têm valor ajustado igual a  $\hat{y}_{1j} = \bar{y}_1 = 59.14$ . O resíduo da primeira observação do primeiro grupo será  $e_{11} = 62.6 - 59.14 = 3.46$  e o da segunda observação desse grupo é  $e_{12} = 59.6 - 59.14 = 0.46$ . De forma análoga, os valores ajustados de qualquer observação no segundo grupo são dados por  $\hat{y}_{2j} = \bar{y}_2 = 46.04$ . O resíduo da terceira observação do segundo grupo é assim  $e_{23} = y_{23} - \bar{y}_2 = 47.5 - 46.04 = 1.46$ . Para calcular a totalidade dos resíduos podemos recorrer ao R (arredondando a três casas decimais):

```
> round(residuals(C02.aov), d=3)
 1 2 3 4 5 6 7 8 9 10 11 12 13
3.46 0.46 5.36 0.16 -0.54 5.46 -8.24 -2.94 -6.84 3.66 4.86 -1.74 1.46
14 15 16 17 18 19 20 21 22 23 24 25 26
3.46 2.46 4.36 -10.84 3.86 -3.44 -4.44 9.05 4.65 -6.65 1.85 3.75 2.05
27 28 29 30 31 32 33 34 35 36 37 38 39
-6.25 -9.45 3.55 -2.55 4.03 -2.67 -6.27 -4.87 3.73 -1.37 -2.87 7.23 -1.07
40 41 42 43 44 45 46 47 48 49 50
4.13 8.46 0.76 -8.64 -5.94 1.36 5.66 6.16 0.36 -0.54 -7.64
```

Com o auxílio do R, podemos obter os dois gráficos de resíduos já considerados no estudo dos modelos de Regressão Linear, através do comando:

```
> plot(C02.aov, which=c(1,2), pch=16, col="blue")
```





O gráfico da esquerda é o gráfico de resíduos usuais (no eixo vertical) vs. valores ajustados da variável resposta (eixo horizontal). O facto de os resíduos surgirem “empilhados” em colunas resulta do já referido facto de todas as observações dum dado nível terem o mesmo valor ajustado  $\hat{y}_{ij} = \bar{y}_i$ , logo, a mesma coordenada no eixo horizontal. Neste caso, observam-se  $k = 5$  colunas. Não parece existir problema com a hipótese de homogeneidade das variâncias, uma vez que a variabilidade dos resíduos não parece diferir muito nos cinco níveis do factor. O *qq-plot* (gráfico à direita) não indicia problemas graves com a Normalidade, dada a disposição aproximadamente linear dos pontos.

Os restantes diagnósticos que foram considerados aquando do estudo da regressão (distâncias de Cook, efeito alavanca) são geralmente de menor utilidade no contexto duma ANOVA. Para as distâncias de Cook sabe-se de antemão qual o efeito de retirar uma observação: além de desequilibrar um delineamento equilibrado, afectará a média das observações no mesmo nível do factor (ou seja, os valores ajustados  $\hat{y}$  nesse nível). Assim valores elevados da distância de Cook correspondem a observações atípicas (*outliers*) no seio dum dado nível. Mas para identificar tais observações, basta o gráfico usual de resíduos contra  $\hat{y}$ . E é possível mostrar que o efeito alavanca de qualquer observação  $y_{ij}$  numa ANOVA a um factor é dada por  $\frac{1}{n_i}$ , onde  $n_i$  indica o número de observações no nível  $i$  da observação. Em delineamentos equilibrados, esse valor é igual para todas as observações (no nosso caso, todas teriam efeito alavanca igual a  $\frac{1}{10}$ ). O gráfico obtido no R com a opção `which=5` é diferente do obtido numa regressão linear: no eixo horizontal indicam-se apenas os diferentes níveis do factor (ordenados por ordem crescente das médias  $\bar{y}_i$ ), uma vez que um gráfico análogo ao construído na regressão linear apenas empilharia todos os resíduos numa única coluna.

- (b) Nesta alínea pede-se para aproveitar os valores das concentrações de  $CO_2$  utilizadas, e tratar essa variável preditora como uma variável numérica, estudando a regressão linear simples de `variacao.massa` sobre `C02.numerico`. O gráfico pedido pode ser construído com o seguinte comando do R. O resultado é mostrado na alínea seguinte.

```
> plot(variacao.massa ~ C02.numerico, data=C02, pch=16)
```

A regressão linear correspondente é dada por:

```
> C02.lm <- lm(variacao.massa ~ C02.numerico, data=C02)
> summary(C02.lm)
```

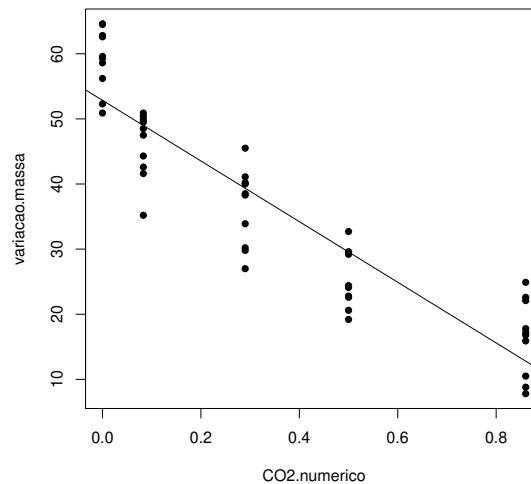
Coefficients:

|              | Estimate | Std. Error | t value | Pr(> t )   |
|--------------|----------|------------|---------|------------|
| (Intercept)  | 52.849   | 1.408      | 37.52   | <2e-16 *** |
| C02.numerico | -46.569  | 3.030      | -15.37  | <2e-16 *** |

---

Residual standard error: 6.637 on 48 degrees of freedom  
 Multiple R-squared: 0.8312, Adjusted R-squared: 0.8276  
 F-statistic: 236.3 on 1 and 48 DF, p-value: < 2.2e-16

A nuvem de pontos pedida, já com a recta de regressão (traçada com o comando `abline(CO2.lm)`) é:



Apesar de alguma tendência para uma relação curvilínea, uma regressão linear simples pode constituir uma modelação aproximada da relação entre concentrações de dióxido de carbono e variação na massa das culturas de *Pseudomonas fragi* (repare-se como seria impossível tirar esta relação se o número de níveis fosse mais pequeno, *e.g.*,  $k = 3$ ). O valor do coeficiente de determinação é claramente significativo ( $p < 2.2 \times 10^{-16}$ ) e bastante elevado ( $R^2 = 0.8312$ ), explicando mais de 83% da variabilidade total observada na variável resposta. Os testes  $F$  de ajustamento global do contexto regressão linear simples e do contexto ANOVA a um factor, não são os mesmos. Como se viu nas aulas teóricas, a ANOVA a um factor pode ser vista como um modelo linear, mas em que as variáveis preditoras são as indicatrizes dos níveis (excepto o primeiro) do factor. Assim, a informação disponível para prever os valores da variável resposta é, no caso da regressão considerada nesta alínea, a variável `CO2.numerico`, com valores numéricos diferentes em cada nível (embora repetidos para as observações dum mesmo nível). No caso da ANOVA a um factor, é o conjunto das indicatrizes de nível e o vector dos  $n$  uns. Sendo diferente a informação preditora, serão diferentes os valores ajustados e os valores dos respectivos  $F_{calc}$  e coeficientes de determinação. Em relação a este último, e embora não seja hábito utilizá-lo no contexto duma ANOVA a um factor, o seu valor é aqui  $R^2 = 0.9003$ , superior ao que se obteve na regressão ( $R^2 = 0.8312$ ), como se pode constatar através do ajustamento obtido utilizando simultaneamente o comando `lm` e o factor preditor `CO2.factor`:

```
> summary(lm(variacao.massa ~ CO2.factor, data=C02))
(...)
Residual standard error: 5.266 on 45 degrees of freedom
Multiple R-squared: 0.9003, Adjusted R-squared: 0.8915
F-statistic: 101.6 on 4 and 45 DF, p-value: < 2.2e-16
```

Repare-se como o valor da estatística calculada,  $F_{calc} = 101.6$ , é o que foi obtido usando o comando `aov`. Um comentário final: o modelo ANOVA não permite, ao contrário da regressão, fazer previsões sobre as variações de massa com concentrações de  $CO_2$  não observadas

na experiência, uma vez que os níveis do factor  $CO_2$  não têm escala (são apenas categorias diferentes).

4. Trata-se dum delineamento factorial a dois factores (**terreno** e **variedade**), mas com uma única observação em cada célula (em cada terreno, apenas há uma parcela com cada variedade). Logo, só é possível ajustar um modelo a dois factores sem interacção.

- (a) A tabela-resumo correspondente é:

```
> terrenos.aov <- aov(rend ~ variedade + terreno, data=terrenos)
> summary(terrenos.aov)
```

|           | Df | Sum Sq | Mean Sq | F value | Pr(>F)     |
|-----------|----|--------|---------|---------|------------|
| variedade | 3  | 1.799  | 0.5997  | 6.145   | 0.00175 ** |
| terreno   | 12 | 2.407  | 0.2006  | 2.056   | 0.04737 *  |
| Residuals | 36 | 3.513  | 0.0976  |         |            |

Desta tabela depreende-se que, aos níveis de significância usuais, deve considerar-se a existência de efeitos do factor variedade:

**Hipóteses:**  $H_0 : \alpha_i = 0, \forall i = 2, 3, 4$  vs.  $H_1 : \exists i = 2, 3, 4$  tal que  $\alpha_i \neq 0$ .

**Estatística do teste:**  $F = \frac{QMA}{QMRE} \cap F_{(a-1, n-(a+b-1))}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica (Unilateral Direita):** Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(3,36)} \approx 2.87$ .

**Conclusões:**  $F_{calc} = 6.145$ , um valor significativo mesmo ao nível  $\alpha = 0.005$ . Logo, rejeita-se  $H_0$  a favor da hipótese de que existem efeitos do factor. Assim, é de concluir que diferentes variedades estejam associadas a diferentes rendimentos médios.

- (b) Um teste aos efeitos do factor **terreno** permite tirar a conclusão que os efeitos deste factor são menos importantes que os efeitos do factor **variedade**, embora ao nível de significância  $\alpha = 0.05$  sejam (por pouco) significativos. Assim,

**Hipóteses:**  $H_0 : \beta_j = 0, \forall j = 2, \dots, 13$  vs.  $H_1 : \exists j = 2, \dots, 13$  tal que  $\beta_j \neq 0$ .

**Estatística do teste:**  $F = \frac{QMB}{QMRE} \cap F_{(b-1, n-(a+b-1))}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica (Unilateral Direita):** Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(12,36)} \approx 2.04$ .

**Conclusões:**  $F_{calc} = 2.056$ , um valor significativo (por muito pouco) ao nível  $\alpha = 0.05$ . Logo, rejeita-se  $H_0$  a favor da hipótese de que existem efeitos do factor **terreno**.

**NOTA:** Num caso como este, em que a conclusão depende do nível de significância usado, é especialmente importante que eventuais fontes de variabilidade, exteriores ao factor sob estudo, mas que afectem a variável resposta, sejam tidas em conta, de forma a reduzir a variabilidade não explicada pelo modelo, isto é, o valor de  $QMRE$ .

5. (a) Trata-se dum delineamento factorial a dois factores, sendo a variável resposta  $Y$  a altura aos dois anos (em cm) dos pinheiros; o primeiro factor (A) a proveniência, com  $a = 5$  níveis e o segundo factor (B) o local do ensaio (com  $b = 2$  níveis). O delineamento é equilibrado, uma vez que em cada uma das  $ab = 10$  células (situações experimentais) existem  $n_c = 6$  observações, num total de  $n = n_c ab = 60$  observações. Existem repetições nas células, logo é possível (e desejável) estudar a existência de eventuais efeitos de interacção.

O modelo ajustado é o modelo ANOVA a dois factores, com efeitos de interacção. Admite-se que os níveis de cada factor estão ordenados por ordem alfabética (que corresponde à ordem em que aparecem no enunciado). Eis o modelo:

- $Y_{ijk} = \mu_{11} + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}$ , para qualquer  $i = 1, 2, 3, 4, 5$ ,  $j = 1, 2$  e  $k = 1, 2, 3, 4, 5, 6$ , sendo  $\mu_{11}$  a altura esperada (aos dois anos) dos pinheiros gregos em Sines;  $\alpha_i$  o efeito principal (acréscimo à altura) associado à proveniência  $i$  (com a restrição  $\alpha_1 = 0$ );  $\beta_j$  o efeito principal (acréscimo à altura) associado a  $j = 2$  (dada a restrição  $\beta_1 = 0$ );  $(\alpha\beta)_{ij}$  o efeito de interacção, isto é, o acréscimo na altura específico da combinação da proveniência  $i$  com o local  $j$ . Dadas as restrições  $(\alpha\beta)_{ij} = 0$  se  $i = 1$  e/ou  $j = 1$ , o modelo apenas prevê efeitos de interacção nas situações experimentais correspondentes a Tavira ( $j = 2$ ) e para proveniências diferentes da Grécia ( $i > 1$ ). Finalmente  $\epsilon_{ijk}$  é o erro aleatório da observação  $Y_{ijk}$ .
  - Admite-se que os erros aleatórios são todos Normais, de média zero e variâncias homogéneas:  $\epsilon_{ijk} \cap \mathcal{N}(0, \sigma^2)$ , para qualquer  $i, j, k$ .
  - Admite-se que os erros aleatórios  $\epsilon_{ijk}$  são independentes.
- (b) Tratando-se dum modelo ANOVA factorial, a dois factores com interacção, a tabela-resumo terá de ter quatro linhas, correspondentes aos três tipos de efeitos previstos (principal de cada factor e de interacção), bem como à variabilidade residual e, opcionalmente, uma quinta linha associada à variabilidade total. A tabela terá as habituais colunas de graus de liberdade, Somas de Quadrados, Quadrados Médios e valor das estatísticas  $F$ . Vejamos como se pode preencher esta tabela.

Sabemos que, neste tipo de modelo, os graus de liberdade associados a  $QMRE$  são dados por  $n - ab$ , onde  $n = 60$  é o número total de observações e  $ab = 10$  é o número de parâmetros existentes no modelo. Assim,  $g.l.(SQRE) = 50$ . Sabemos ainda que, para os vários tipos de efeitos, os graus de liberdade são dados pelo número de parcelas de cada tipo de efeito, após a introdução das restrições, ou seja, associado a  $SQA$  há  $a - 1 = 4$  g.l., associado a  $SQB$  há  $b - 1 = 1$  g.l., e associado a  $SQAB$  há  $(a - 1)(b - 1) = 4$  graus de liberdade.

No enunciado é dada a Soma de Quadrados associada ao que foi designado factor A, tendo-se  $SQA = 280.61$ , donde se conclui que  $QMA = \frac{SQA}{a-1} = \frac{280.61}{4} = 70.1525$ . No enunciado é também dado o Quadrado Médio Residual, tendo-se  $QMRE = 16.59$ , donde  $SQRE = QMRE \times (n - ab) = 16.59 \times 50 = 829.50$ . Ora, sabemos pelo formulário que:

$$\begin{aligned} SQB &= a n_c \sum_{j=1}^2 (\bar{y}_{.j} - \bar{y}_{...})^2 \\ &= 5 \times 6 \times [(28.14 - 31.76298)^2 + (35.38 - 31.76298)^2] = 786.2645 . \end{aligned}$$

Donde  $QMB = \frac{SQB}{b-1} = 786.2645$ . O enunciado refere ainda a variância da totalidade das 60 observações,  $s_y^2 = 34.49584$ , donde se pode concluir que a Soma de Quadrados Total é  $SQT = (n - 1) s_y^2 = 59 \times 34.49584 = 2035.255$ . Uma vez que sabemos que esta Soma de Quadrados Total se pode decompor como  $SQT = SQA + SQB + SQAB + SQRE$ , torna-se possível calcular  $SQAB = SQT - (SQA + SQB + SQRE) = 2035.255 - (280.61 + 786.2645 + 829.50) = 138.8801$ . Assim, o Quadrado Médio associado à interacção é dado por  $QMAB = \frac{SQAB}{(a-1)(b-1)} = \frac{138.8801}{4} = 34.7200$ .

Finalmente, os valores das estatísticas  $F$  são dados, para os três tipos de efeitos, pela razão entre o Quadrado Médio do referido tipo de efeito e  $QMRE$ . A tabela completa fica assim:

|              | g.l. | Soma de Quadrados | Quadrado Médio | F      |
|--------------|------|-------------------|----------------|--------|
| Proveniência | 4    | 280.61            | 70.1525        | 4.229  |
| Local        | 1    | 786.2645          | 786.2645       | 47.394 |
| Interacção   | 4    | 138.8801          | 34.7200        | 2.093  |
| Residual     | 50   | 829.50            | 16.59          | –      |

- (c) Vai-se efectuar em pormenor o teste aos efeitos principais do Factor A (proveniência dos pinheiros), e descrever sinteticamente os testes aos efeitos principais do Factor B (local) e aos efeitos de interacção.

**Hipóteses:**  $H_0 : \alpha_i = 0, \forall i$  vs.  $H_1 : \exists i$  tal que  $\alpha_i \neq 0$ .

**Estatística do Teste:**  $F_A = \frac{QMA}{QMRE} \cap F_{[a-1, n-ab]}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica:** (Unilateral direita) Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(4,50)} \approx 2.57$  (entre os valores tabelados 2.53 e 2.61).

**Conclusões:** Como  $F_{calc} = \frac{QMA}{QMRE} = 4.229 > 2.57$ , rejeita-se  $H_0$ , sendo possível concluir pela existência de efeitos principais de proveniência (ao nível  $\alpha = 0.05$ ).

No teste aos efeitos principais do factor local do estudo, as hipóteses do teste podem ser escritas apenas como  $H_0 : \beta_2 = 0$  vs.  $H_1 : \beta_2 \neq 0$ , uma vez que após a imposição da restrição  $\beta_1 = 0$ , apenas sobra um efeito deste tipo, o efeito  $\beta_2$  associado a Tavira. O valor calculado da estatística de teste é muito grande ( $F_{calc} = 47.394$ ) deixando antever a rejeição de  $H_0$ , facto que é confirmado determinando nas tabelas o limiar da região crítica unilateral direita:  $f_{0.05(1,50)} \approx 4.04$  (entre os valores tabelados 4.00 e 4.08). Assim, conclui-se claramente pela existência de efeitos principais de localidade, o que neste caso significa que existe um efeito associado à passagem do local de plantação de Sines para Tavira. Uma rápida inspecção das médias de local sugere que se trata dum maior crescimento dos pinheiros em Tavira, pelo que se deduz que  $\beta_2$  terá um valor positivo.

No teste aos efeitos de interacção, com hipóteses  $H_0 : (\alpha\beta)_{ij} = 0$ , para todo o  $i$  e  $j$ , contra a hipótese alternativa de que existe pelo menos uma célula  $(i, j)$  onde  $(\alpha\beta)_{ij} \neq 0$ , o valor calculado da estatística de teste é  $F_{calc} = 2.093$ , inferior ao limiar da região crítica, que é (por coincidência) igual ao do teste aos efeitos do factor A,  $f_{0.05(4,50)} \approx 2.57$ . Logo, não se rejeita  $H_0$  (para  $\alpha = 0.05$ ), e conclui-se pela inexistência de efeitos significativos de interacção.

- (d) **[Material Complementar]** Nesta alínea é pedido para verificar se o facto da maior altura média amostral de Sines (31.16, para pinheiros provenientes de Marrocos) ser menor que a mais baixa altura média amostral em Tavira (33.56, para pinheiros da segunda proveniência italiana) é uma relação que se possa estender à população. Vamos responder efectuando, como solicitado no enunciado, um teste de Tukey, e usando  $\alpha = 0.05$ . Ora, o termo de comparação é (como indicado no formulário e usando as tabelas da distribuição de Tukey):

$$q_{\alpha(ab, n-ab)} \sqrt{\frac{QMRE}{n_c}} = q_{0.05(10,50)} \sqrt{\frac{16.59}{6}} = 4.68 \times 1.662829 = 7.782039 .$$

Ora, a diferença entre as médias amostrais das duas células referidas acima é apenas  $|31.16 - 33.56| = 2.40$ , logo inferior ao termo de comparação, pelo que não é uma diferença significativa (ao nível  $\alpha = 0.05$ ). Assim, não é possível afirmar que as médias populacionais em Tavira sejam sempre maiores às de Sines, independentemente das proveniências. Alguns pares de médias populacionais podem ser consideradas diferentes (por exemplo, o crescimento médio dos pinheiros gregos em Sines e em Tavira), mas será preciso levar em conta as proveniências, e não apenas o local da realização do estudo.

6. (a) Trata-se dum delineamento factorial a dois factores: *localidade* (Factor A, com  $a = 4$  níveis) e *cultivar* (Factor B, com  $b = 9$  níveis). Existem  $n_{ij} = 4 = n_c$  repetições em todas as  $ab = 36$  situações experimentais (células), pelo que se trata dum delineamento equilibrado. Existem ao todo  $n = abn_c = 144$  observações da variável resposta  $Y$  (rendimento, em  $kg/ha$ ). O modelo ANOVA adequado é o modelo ANOVA a dois factores, com interacção, dado por:
- i.  $Y_{ijk} = \mu_{11} + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}$ ,  $\forall i = 1, 2, 3, 4$ ,  $j = 1, 2, \dots, 9$ ,  $k = 1, 2, 3, 4$ , com  $\alpha_1 = 0$ ,  $\beta_1 = 0$ ,  $(\alpha\beta)_{1j} = 0$  para qualquer  $j$ , e  $(\alpha\beta)_{i1} = 0$  para qualquer  $i$ , onde
    - $Y_{ijk}$  indica o rendimento na  $k$ -ésima parcela da localidade  $i$ , associada à cultivar  $j$ ;
    - $\mu_{11}$  indica o rendimento médio (populacional) da cultivar *Celta*, em Elvas;
    - $\alpha_i$  indica o efeito principal da localidade  $i$ ;
    - $\beta_j$  indica o efeito principal da cultivar  $j$ ;
    - $(\alpha\beta)_{ij}$  indica o efeito de interacção entre a localidade  $i$  e a cultivar  $j$ ; e
    - $\epsilon_{ijk}$  indica o erro aleatório associado à observação  $Y_{ijk}$ .
  - ii.  $\epsilon_{ijk} \cap \mathcal{N}(0, \sigma^2)$ ,  $\forall i, j, k$ .
  - iii.  $\{\epsilon_{ijk}\}_{i,j,k}$  constitui um conjunto de variáveis aleatórias independentes.
- (b) i. Os nove valores em falta na tabela são dados por:
- $g.l.(SQA) = a - 1 = 3$ ;
  - $g.l.(SQB) = b - 1 = 8$ ;
  - $g.l.(SQAB) = (a - 1)(b - 1) = 3 \times 8 = 24$ ;
  - $g.l.(SQRE) = n - ab = 144 - 36 = 108$ ;
  - $SQB = QMB(b - 1) = 964\,060 \times 8 = 7\,712\,480$ ;
  - $SQAB = SQT - (SQA + SQB + SQRE) = (n - 1)s_y^2 - 219\,628\,472 = 143 \times 1\,714\,242 - 219\,628\,472 = 25\,508\,134$ ;
  - $QMA = \frac{SQA}{a-1} = \frac{183\,759\,916}{3} = 61\,253\,305$ ;
  - $QMAB = \frac{SQAB}{(a-1)(b-1)} = \frac{25\,508\,134}{24} = 1\,062\,839$ ;
  - $F_B = \frac{QMB}{QMRE} = \frac{964\,060}{260\,704} = 3.69791$ .
- ii. Em qualquer modelo linear (regressão ou ANOVA), a variância dos erros aleatórios do modelo ( $V[\epsilon_i] = \sigma^2$ ) é estimado pelo Quadrado Médio Residual. No nosso caso, a estimativa de  $\sigma^2$  é dada no enunciado:  $QMRE = 260\,704$ . O valor muito elevado nada indica de especial, uma vez que a sua interpretação tem de levar em conta as unidades de medida dos dados, que são  $(kg\ ha^{-1})^2$ . De facto sabemos pelo enunciado que as unidades de medida da variável resposta são  $kg/ha$ . Sabemos que os resíduos ( $e_i = y_i - \hat{y}_i$ ) têm as mesmas unidades de medida que a variável resposta. Sabemos que o QMRE é a Soma de Quadrados dos Resíduos a dividir pelos graus de liberdade associados, pelo que as unidades de medida do QMRE são o quadrado das unidades de medida da variável resposta. Bastava que os valores da variável resposta tivessem sido medidos em toneladas por hectare, para que o Quadrado Médio Residual viesse em  $(t\ ha^{-1})^2$ , ou seja, que fosse um milhão de vezes inferior ao valor acima indicado:  $QMRE = 0.260704$ . Mas isso não altera os dados, nem a significância de cada tipo de efeitos previsto no modelo. Assim, não é possível avaliar a estimativa de  $\sigma^2$  apenas olhando para o valor absoluto de  $QMRE$ : é essencial ter em conta as unidades de medida associadas.
  - iii. Pedem-se os três testes  $F$  para cada tipo de efeitos previstos no modelo. Efectuemos em pormenor o teste à existência de efeitos de interacção entre localidade e cultivar:

**Hipóteses:**  $H_0 : (\alpha\beta)_{ij} = 0, \forall i = 2, 3, 4$  e  $j = 2, 3, \dots, 9$  [não há interacção]  
vs.  $H_1 : \exists i = 2, 3, 4, j = 2, 3, \dots, 9$  tais que  $(\alpha\beta)_{ij} \neq 0$  [há interacção].

**Estatística do teste:**  $F = \frac{QMAB}{QMRRE} \cap F_{[(a-1)(b-1), n-ab]}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.01$ .

**Região Crítica (Unilateral Direita):** Rejeitar  $H_0$  se  $F_{calc} > f_{0.01(24,108)} \approx 1.97$ .

**Conclusões:** O valor da estatística do teste foi calculado na alínea anterior:  $F_{calc} = 4.0768$ . É um valor significativo ao nível  $\alpha = 0.01$ , rejeitando-se  $H_0$  a favor da hipótese alternativa de que existem efeitos de interacção entre localidade e cultivar.

No que respeita ao teste para os efeitos principais do factor *localidade*, as hipóteses em confronto são  $H_0 : \alpha_i = 0, \forall i = 2, 3, 4$  vs.  $H_1 : \exists i = 2, 3, 4$ , tal que  $\alpha_i \neq 0$ . A Região Crítica é agora dada pela rejeição de  $H_0$  caso  $F_{calc} > f_{0.01(3,108)} \approx 3.97$ . O valor elevadíssimo da estatística calculada  $F_{calc} = 234.9531$  leva à rejeição clara de  $H_0$ , concluindo-se pela existência de importantes efeitos de localidade, nos rendimentos.

Finalmente, no teste aos efeitos principais do factor *cultivar*, as hipóteses em confronto são  $H_0 : \beta_j = 0, \forall j = 2, 3, \dots, 9$  vs.  $H_1 : \exists j = 2, 3, \dots, 9$ , tal que  $\beta_j \neq 0$ . A Região Crítica é agora dada pela rejeição de  $H_0$  caso  $F_{calc} > f_{0.01(8,108)} \approx 2.68$ . O valor da estatística calculada  $F_{calc} = 3.698$  pertence à Região Crítica, levando à rejeição de  $H_0$ , concluindo-se também pela existência de efeitos de cultivar sobre os rendimentos.

Assim, conclui-se pela existência dos três tipos de efeitos, ao nível  $\alpha = 0.01$ , com destaque para a existência clara de efeitos de localidade.

- iv. Pedese para discutir o efeito sobre a tabela resultante de dividir a variável resposta por mil (passando o rendimento a ser expresso em  $t/ha$ ). Os graus de liberdade não são, naturalmente, afectados. O mesmo não se passa com as Somas de Quadrados. À nova variável  $Y^* = Y/1000$  corresponderão novas médias de nível, de célula e global, que também resultam de dividir por mil (para ficarem em  $t/ha$ ). Tendo em conta que no modelo em questão, as médias de célula definem os valores ajustados, tem-se  $\hat{Y}_{ijk}^* = \hat{Y}_{ijk}/1000$ . Assim, as novas Somas de Quadrados resultam de dividir as suas congéneres originais por  $1000^2$ , ou seja, por um milhão. De facto,  $SQT^* = \sum_i \sum_j \sum_k (Y_{ijk}^* - \bar{Y}_{...}^*)^2 = \sum_i \sum_j \sum_k (Y_{ijk}/1000 - \bar{Y}_{...}/1000)^2 = SQT/(1000^2)$ . Também  $SQRE^* = \sum_i \sum_j \sum_k (Y_{ijk}^* - \hat{Y}_{ijk}^*)^2 = \sum_i \sum_j \sum_k (Y_{ijk}/1000 - \hat{Y}_{ijk}/1000)^2 = SQRE/(1000^2)$ . De forma análoga, e utilizando as fórmulas para delineamentos equilibrados,

$$SQA^* = bn_c \sum_{i=1}^a (\bar{Y}_{i..}^* - \bar{Y}_{...}^*)^2 = bn_c \sum_{i=1}^a (\bar{Y}_{i..}/1000 - \bar{Y}_{...}/1000)^2 = SQA/(1000^2)$$

$$SQB^* = an_c \sum_{j=1}^b (\bar{Y}_{.j.}^* - \bar{Y}_{...}^*)^2 = an_c \sum_{j=1}^b (\bar{Y}_{.j.}/1000 - \bar{Y}_{...}/1000)^2 = SQB/(1000^2).$$

Por diferença, tem igualmente de verificar-se  $SQAB^* = SQAB/(1000^2)$ . Assim, toda a coluna de Somas de Quadrados na tabela será dividida por um milhão. Essa mesma transformação aplica-se à coluna de Quadrados Médios (que resulta de dividir Somas de Quadrados por graus de liberdade). Mas na coluna final, correspondente aos valores calculados das estatísticas  $F$ , o quociente de Quadrados Médios mantém-se inalterado (a transformação multiplicativa de numerador e denominador é igual). Logo, as conclusões de todos os testes (incluindo os respectivos *p-values*) mantêm-se inalterados.

v. **[Material Complementar]** Os dois gráficos de interacção reflectem a mesma informação, embora de formas diferentes. No gráfico da esquerda, as quatro localidades definem posições no eixo horizontal. Por cima de cada localidade encontram-se nove pontos, associados às nove cultivares. A ordenada de cada um desses nove pontos é dada pelo rendimento médio das parcelas correspondentes a essa combinação de localidade e cultivar. Os segmentos de recta unem os pontos correspondentes a cada cultivar (segundo a legenda indicada no gráfico). Embora haja algum paralelismo nas nove curvas seccionalmente lineares, para as três primeiras localidades, os rendimentos na Revilheira sugerem a existência de efeitos de interacção. Por exemplo, a cultivar *TE9110*, que regista o rendimento mais baixo em Elvas (facto que se pode confirmar na tabela de médias dada na alínea c) tem o segundo mais elevado rendimento na Revilheira. Também a cultivar *Celta*, cujo rendimento em Benavila é o terceiro mais baixo, regista o segundo maior rendimento em Elvas. Assim, há cultivares que manifestam “preferências” ou “aversões” por diferentes localidades, reflectindo efeitos de interacção. O teste à interacção efectuado na alínea anterior confirma que esses efeitos são significativos, ao nível  $\alpha = 0.01$ .

O gráfico da direita dá, como se disse, uma perspectiva diferente sobre a mesma informação. Agora, são as cultivares que definem nove posições no eixo horizontal. Por cima de cada uma dessas posições (cultivares) há quatro pontos, com ordenadas dadas pelos rendimentos médios da referida cultivar, nas quatro localidades consideradas no ensaio. Segmentos de recta unem os pontos correspondentes a uma mesma localidade. Neste gráfico torna-se evidente que os rendimentos são sempre bastante superiores em Elvas (no gráfico da esquerda, esse facto reflectia-se no “pico” por cima de Elvas). Essa será a principal razão pela clara rejeição da hipótese nula no teste à existência de efeitos principais de localidade. Por outro lado, os efeitos de interacção reflectem-se na mais visível ausência de paralelismo, nomeadamente nos traços correspondentes a Elvas e Revilheira, que para várias cultivares parecem ter comportamentos quase antagónicos.

7. (a) Trata-se dum delineamento factorial a dois factores: *Temperatura de conservação* (Factor A), com  $a = 2$  níveis, e *Tempo de armazenamento* (Factor B), com  $b = 4$  níveis. Para modelar a variável resposta  $Y$  (alterações no conteúdo em taninos das polpas de sapoti), utiliza-se um modelo ANOVA a dois factores, com interacção. É possível estudar a interacção devido à presença de repetições nas  $2 \times 4 = 8$  células. Sempre que possível, é desejável considerar este modelo para delineamentos factoriais a dois factores, deixando que sejam os dados a sugerir se se deve admitir a existência desse tipo de efeitos. O delineamento é equilibrado, uma vez que todas as células têm o mesmo número de repetições:  $n_{ij} = 4 = n_c$  ( $\forall i, j$ ), para um total de  $n = 8 \times 4 = 32$  observações. O modelo é dado por:

- i.  $Y_{ijk} = \mu_{11} + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}$ ,  $\forall i = 1, 2$ ,  $j = 1, 2, 3, 4$ ,  $k = 1, 2, 3, 4$ ,  
com  $\alpha_1 = 0$ ,  $\beta_1 = 0$ ,  $(\alpha\beta)_{1j} = 0$  para qualquer  $j$ , e  $(\alpha\beta)_{i1} = 0$  para qualquer  $i$ , onde
- $Y_{ijk}$  indica a  $k$ -ésima observação (repetição) na célula definida pelo nível  $i$  do Factor A e o nível  $j$  do Factor B;
  - $\mu_{11}$  indica a média (populacional) das observações na célula (1,1), ou seja, com temperatura alta e 0 dias de armazenamento;
  - $\alpha_i$  indica o efeito do nível  $i$  do Factor A (*Temperatura*);
  - $\beta_j$  indica o efeito do nível  $j$  do Factor B (*Tempo de armazenamento*);
  - $(\alpha\beta)_{ij}$  indica o efeito de interacção na célula  $(i, j)$ ; e



- $\epsilon_{ijk}$  indica o erro aleatório associado à observação  $Y_{ijk}$ .
  - ii.  $\epsilon_{ijk} \cap \mathcal{N}(0, \sigma^2)$ ,  $\forall i, j, k$ .
  - iii.  $\{\epsilon_{ijk}\}_{i,j,k}$  constituem um conjunto de variáveis aleatórias independentes.
- (b) A tabela-resumo desta ANOVA terá três linhas associadas a cada tipo de efeitos previsto no modelo (ou seja, efeitos principais do Factor A, efeitos principais do Factor B e efeitos de interacção) e ainda uma linha para o residual (podendo também incluir-se a linha associada à variabilidade Total). Como em qualquer modelo ANOVA, a tabela-resumo tem as seguintes colunas: Somas de Quadrados, graus de liberdade correspondentes, Quadrados Médios e estatísticas  $F$ . Os graus de liberdade são dados por:
- Factor A:  $a - 1 = 1$ ;
  - Factor B:  $b - 1 = 3$ ;
  - Interacção:  $(a - 1)(b - 1) = 3$ ;
  - Residual:  $n - ab = 32 - 8 = 24$ .

Para calcular as Somas de Quadrados, registamos que no enunciado é dada a Soma de Quadrados Residual  $SQRE = 20.72$ . É igualmente dado o Quadrado Médio do Factor B, e multiplicando pelos respectivos graus de liberdade obtém-se  $SQB = QMB(b - 1) = 96.01 \times 3 = 288.03$ . A Soma de Quadrados Total também pode ser calculada facilmente, uma vez que no enunciado é dada a variância da totalidade das observações de  $Y$ ,  $s_y^2 = 47.83222$ , e  $SQT = (n - 1)s_y^2 = 31 \times 47.83222 = 1482.799$ . Assim, faltam as duas Somas de Quadrados relativas aos efeitos principais do factor A ( $SQA$ ) e aos efeitos de interacção ( $SQAB$ ). Utilizando a expressão para  $SQA$ , no caso de delineamentos equilibrados (disponível no formulário) e os valores das médias de nível do factor A e da média geral (disponíveis no enunciado), tem-se  $SQA = bn_c \sum_{i=1}^a (\bar{Y}_{i..} - \bar{Y}_{...})^2 = 16 [(24.681 - 22.14375)^2 + (19.606 - 22.14375)^2] = 16 \times 12.87781 = 206.045$ . A última Soma de Quadrados em falta ( $SQAB$ ) pode ser calculada a partir das restantes quatro:  $SQAB = SQT - (SQA + SQB + SQRE) = 1482.799 - (206.045 + 288.03 + 20.72) = 968.004$ . Assim,

| Variacão   | g.l. | SQs      | QMs                                        | $F_{calc}$                         |
|------------|------|----------|--------------------------------------------|------------------------------------|
| Factor A   | 1    | 206.045  | $QMA = \frac{SQA}{a-1} = 206.045$          | $F = \frac{QMA}{QMRE} = 238.6622$  |
| Factor B   | 3    | 288.03   | $QMB = \frac{SQB}{b-1} = 96.01$            | $F = \frac{QMB}{QMRE} = 111.2085$  |
| Interacção | 3    | 968.004  | $QMAB = \frac{SQAB}{(a-1)(b-1)} = 322.668$ | $F = \frac{QMAB}{QMRE} = 373.7467$ |
| Residual   | 24   | 20.72    | $QMRE = \frac{SQRE}{n-ab} = 0.8633333$     | -                                  |
| Total      | 31   | 1482.799 | -                                          | -                                  |

- (c) De acordo com o modelo, a influência do Factor B nos valores da variável resposta pode resultar de dois tipos de efeitos: os efeitos principais do Factor B (os  $\beta_j$ ) ou os efeitos de interacção (os  $(\alpha\beta)_{ij}$ ). Efectuaremos estes dois testes, começando pelo dos efeitos de interacção. Neste exemplo, e como o Factor A apenas tem dois níveis, o índice  $i$  nos efeitos de interacção apenas toma o valor  $i = 2$ .

**Hipóteses:**  $H_0 : (\alpha\beta)_{2j} = 0$ ,  $\forall j = 2, 3, 4$  vs.  $H_1 : \exists j = 2, 3, 4$  tal que  $(\alpha\beta)_{2j} \neq 0$ .

**Estatística do teste:**  $F = \frac{QMAB}{QMRE} \cap F_{[(a-1)(b-1), n-ab]}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica (Unilateral Direita):** Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(3,24)} = 3.01$ .

**Conclusões:** O valor da estatística do teste foi calculado na alínea anterior:  $F_{calc} = 373.7467$ . É um valor claramente significativo e rejeita-se  $H_0$  a favor da hipótese alternativa de que existem efeitos de interacção.

Já é possível responder afirmativamente: o Factor B tem efeitos sobre os valores médios de  $Y$ . No entanto, efectuaremos também o teste aos efeitos principais do Factor B:

**Hipóteses:**  $H_0 : \beta_j = 0, \forall j = 2, 3, 4$  vs.  $H_1 : \exists j = 2, 3, 4$  tal que  $\beta_j \neq 0$ .

**Estatística do teste:**  $F = \frac{QMB}{QMRE} \cap F_{(b-1, n-ab)}$ , sob  $H_0$ .

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica (Unilateral Direita):** Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(3,24)} = 3.01$ .

**Conclusões:** O valor da estatística do teste foi calculado na alínea anterior:  $F_{calc} = 111.2085$ . É um valor claramente significativo e rejeita-se  $H_0$  a favor da hipótese de que existem efeitos principais do Factor B.

Assim, quer pela via dos efeitos principais, quer pela via dos efeitos de interacção, o Factor B (*tempo de armazenamento*) afecta os conteúdos médios de taninos nos sapotis.

8. (a) Trata-se dum delineamento a dois factores, o factor *casta* (factor A), e o factor *genótipo* (factor B). O objectivo do estudo é avaliar os eventuais efeitos destes factores sobre a variável resposta (rendimento). Pela própria natureza dos factores em questão, o delineamento deve ser considerado *hierarquizado*, com genótipos subordinados a castas. Não faria sentido considerar o delineamento factorial: não há cruzamentos entre cada um dos oito genótipos e cada uma das duas castas, já que um genótipo apenas faz sentido quando referido à sua casta.

Assim, temos  $a = 2$  castas (níveis do factor A) e, para o factor subordinado genótipos, há  $b_1 = 4$  genótipos para a casta 1 (Antão Vaz) e  $b_2 = 4$  genótipos para a casta 2 (Malvasia Fina). Ao todo há  $b_1 + b_2 = 8$  situações experimentais, e  $n_c = 8$  repetições em cada uma das situações experimentais, num total de  $n = 64$  observações. O modelo mais adequado será o modelo hierarquizado:

- $Y_{ijk} = \mu_{11} + \alpha_i + \beta_{j(i)} + \epsilon_{ijk}, \forall i, j, k$ , onde  $Y_{ijk}$  indica o rendimento da repetição  $k$  ( $k = 1, 2, \dots, 8$ ) do genótipo  $j$  ( $j = 1, 2, 3, 4$ ) da casta  $i$  ( $i = 1, 2$ ). Impõem-se as restrições  $\alpha_1 = 0, \beta_{1(i)} = 0$  para  $i = 1, 2$ . Com estas restrições, o parâmetro  $\mu_{11}$  é o rendimento médio populacional do primeiro genótipo da casta 1, isto é, do genótipo AN105 da casta Antão Vaz;  $\alpha_2$  é o efeito da casta Malvasia Fina;  $\beta_{j(i)}$  ( $j = 2, 3, 4$ ) é o efeito do genótipo  $j$  na casta  $i = 1, 2$ , e  $\epsilon_{ijk}$  é o erro aleatório associado à observação  $Y_{ijk}$ , que corresponde à variabilidade não explicada pelos efeitos previstos no modelo.
- $\epsilon_{ijk} \cap \mathcal{N}(0, \sigma^2)$ , para qualquer  $i, j, k$ .
- Os erros aleatórios  $\epsilon_{ijk}$  são independentes.

- (b) Sabemos que os graus de liberdade na tabela-resumo da ANOVA são dados por:  $a - 1 = 1$  para o efeitos de castas;  $(b_1 - 1) + (b_2 - 1) = 6$  para os efeitos do factor subordinado, genótipos; e  $n - (b_1 + b_2) = 64 - 8 = 56$  para o residual. Por outro lado, conhecemos a partir do enunciado a Soma de Quadrados do Factor A (castas),  $SQA = 79.73597$  e o Quadrado Médio Residual,  $QMRE = \frac{SQRE}{n - (b_1 + b_2)} = 2.873782$ , de onde é possível obter a Soma de Quadrados Residual  $SQRE = 2.873782 \times 56 = 160.9318$ . A Soma de Quadrados associada ao factor subordinado (genótipos) pode ser obtida pela diferença da soma das outras SQs já calculadas em relação à Soma de Quadrados Total, que sai do conhecimento da variância amostral da totalidade das 64 observações. Assim,  $SQT = (n - 1)s_y^2 = 63 \times 5.389415 = 339.5331$ , logo  $SQB(A) = SQT - (SQA + SQRE) = 339.5331 - (79.73597 + 160.9318) = 98.86533$ . Os Quadrados Médios restantes obtêm-se dividindo Somas de Quadrados pelos respectivos graus de liberdade e os valores das duas estatísticas  $F$  resultam de dividir o correspondente quadrado médio pelo  $QMRE$ . Os valores resultantes são sintetizados na tabela em baixo.

| Variaco        | g.l.    | SQs      | QMs       | F                                                 |
|-----------------|---------|----------|-----------|---------------------------------------------------|
| Casta (A)       | 1       | 79.73597 | 79.73597  | $F_A = \frac{79.73597}{2.873782} = 27.74601$      |
| Gentipo [B(A)] | 6       | 98.86533 | 16.47755  | $F_{B(A)} = \frac{16.47755}{2.873782} = 5.733751$ |
| Residual        | 56      | 160.9318 | 2.873782  | –                                                 |
| Total           | 63      | 339.5331 | 5.389415  | –                                                 |
|                 | $(n-1)$ | (SQT)    | $(s_y^2)$ | –                                                 |

- (c) Para responder ser necessrio efectuar um teste  $F$  aos efeitos do factor subordinado (gentipos), cuja hiptese nula corresponde  inexistncia desse tipo de efeitos.

**Hipteses:**  $H_0 : \beta_{j(i)} = 0, \forall i, j$  vs.  $H_1 : \exists i, j$  tal que  $\beta_{j(i)} \neq 0$ .

**Estatística do Teste:**  $F_{B(A)} = \frac{Q_{MB(A)}}{Q_{MRE}} \cap F_{[(b_1-1)+(b_2-1), n-(b_1+b_2)]}$ , sob  $H_0$ .

**Nvel de significncia:** O enunciado pede o nvel  $\alpha = 0.05$ .

**Regio Crtica:** (Unilateral direita) Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(6,56)}$  que, pelas tabelas  um valor entre os valores tabelados 2.25 e 2.34.

**Concluses:** Como  $F_{calc} = 5.733751 > 2.34$ , rejeita-se  $H_0$ , o que corresponde a admitir a existncia de efeitos de gentipos.

Assim, foi importante prever este tipo de efeitos. Ignorar a existncia de efeitos de gentipos iria inflaccionar a Soma de Quadrados Residual, o que poderia mascarar a existncia de efeitos do outro factor (casta), mesmo que eles existam.

- (d) Um teste anlogo, mas aos efeitos do factor dominante (casta) ter como hipteses  $H_0 : \alpha_2 = 0$  (uma vez que apenas existem duas castas e imps-se a restrio  $\alpha_1 = 0$ ) vs.  $H_1 : \alpha_2 \neq 0$ . A regio crtica deste teste (igualmente unilateral direita)   $f_{0.05(1,56)}$ , um valor entre os valores tabelados 4.00 e 4.08. Como  $F_{calc} = 27.746 > 4.08$ , rejeita-se a hiptese nula. Assim, conclui-se (ao nvel de significncia  $\alpha = 0.05$ ) que o efeito  $\alpha_2 \neq 0$ , ou seja que, para alm de existirem efeitos de gentipos, h um efeito significativo de casta, e havendo apenas duas castas, pode-se afirmar que os rendimentos da casta Malvasia Fina so significativamente diferentes dos da casta Anto Vaz.

- (e) **[Material Complementar]** O gentipo MF201 referido no enunciado tem o maior rendimento mdio amostral  $\bar{y}_{2,4} = 7.678$  (ordenando os gentipos como o R). Pretende-se saber que outras mdias amostrais  $\bar{y}_{ij}$  diferem significativamente de  $\bar{y}_{2,4}$ . Utilizaremos as comparaes mltiplas de Tukey ao nvel global  $\alpha = 0.05$ . O termo de comparao correspondente   $q_{\alpha(b_1+b_2, n-(b_1+b_2))} \sqrt{\frac{Q_{MRE}}{n_c}} = q_{0.05(8,56)} \sqrt{\frac{2.873782}{8}} \approx 4.45 \times 0.5993519 = 2.667$ . Qualquer mdia amostral de rendimento de gentipo inferior a  $7.678 - 2.667 = 5.011$  dever assim ser considerada significativamente diferente da mdia do gentipo MF201. H apenas dois gentipos que no tm rendimentos significativamente diferentes, ambos da casta Malvasia Fina: MF1420 e MF1426. Assim, no se rejeitam as hipteses  $\mu_{MF201} = \mu_{MF1420}$  e  $\mu_{MF201} = \mu_{MF1426}$ . Os trs gentipos em questo so da casta Malvasia Fina, o que  coerente com a concluso da alnea anterior: para alm de efeitos de gentipo,  possvel falar de efeitos de casta, sendo os rendimentos da casta Malvasia Fina globalmente superiores.

9. (a) Pede-se para mostrar que a soma dos  $n_i$  resduos  $e_{ij}$ , correspondentes ao nvel  $i$  do Factor ( $i = 1, 2, \dots, k$ ), numa ANOVA a 1 Factor,  nula. Sabemos que, neste tipo de delineamento, os valores ajustados de cada observao correspondem  mdia amostral das  $n_i$  observaes no nvel  $i$  do Factor em que essa observao foi efectuada. Assim,

$$\sum_{j=1}^{n_i} e_{ij} = \sum_{j=1}^{n_i} (y_{ij} - \hat{y}_{ij}) = \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i) = 0,$$

uma vez que se trata duma soma de desvios dum conjunto de observações em relação à sua média (ou seja, do tipo  $\sum_{i=1}^n (x_i - \bar{x})$ , estudada no Exercício 3 da Regressão Linear) que tem sempre soma zero.

- (b) Trata-se duma situação análoga à da alínea anterior. Num modelo ANOVA a dois factores, com efeitos de interacção, sabemos que os valores ajustados  $\hat{y}_{ijk}$  correspondem às médias  $\bar{y}_{ij.}$  das observações da célula da referida observação. Assim, a soma dos resíduos das  $n_{ij}$  observações efectuadas na célula  $(i, j)$  é dada por:

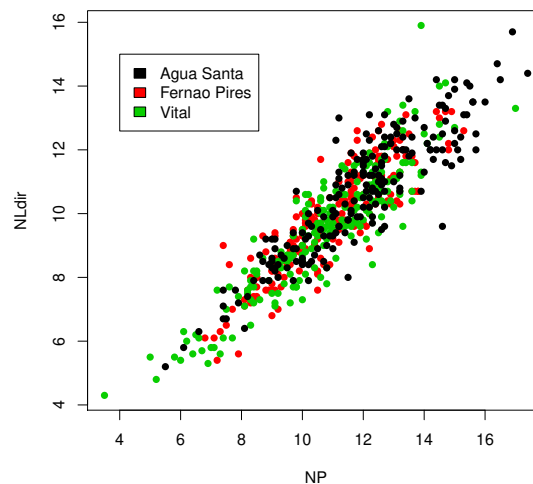
$$\sum_{k=1}^{n_{ij}} e_{ijk} = \sum_{k=1}^{n_{ij}} (y_{ijk} - \hat{y}_{ijk}) = \sum_{k=1}^{n_{ij}} (y_{ijk} - \bar{y}_{ij.}) = 0.$$

### 3 Análise de Covariância

1. Neste exercício consideram-se os dados da *data frame* `videiras`. A variável resposta é, em todas as alíneas, o comprimento da nervura lateral direita (`NLdir`) e o preditor, o comprimento da nervura principal (`NP`).

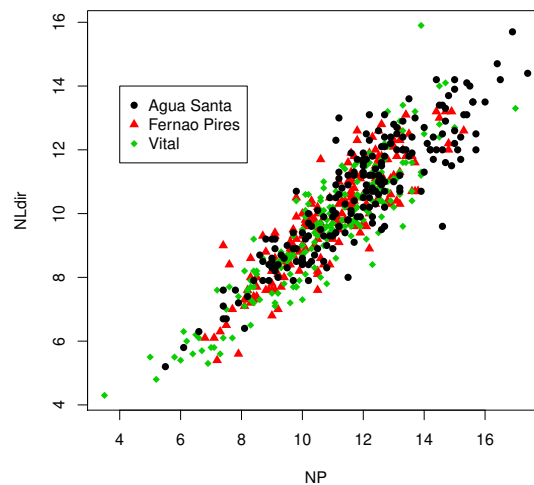
- (a) Os comandos R para obter a nuvem de pontos pedida, e o respectivo resultado, são:

```
> plot(NLdir ~ NP, col=Casta, data=videiras, pch=16)
> legend(4,15,legend=levels(videiras$Casta), fill=1:3)
```



Alternativamente, podemos também querer construir um gráfico com, não apenas cores diferentes, mas também símbolos diferentes para cada casta. Eis uma forma possível de construir um tal gráfico no R, usando os símbolos a que correspondem os códigos 16 (círculos), 17 (triângulos) e 18 (losangos), como indicado na legenda.

```
> plot(NLdir ~ NP, col=as.numeric(Casta), pch=as.numeric(Casta)+15, data=videiras)
> legend(4,14,levels(videiras$Casta),col=1:3, pch=16:18)
```



A nuvem de pontos sugere a existência duma relação linear bastante intensa, que poderá ser a mesma nas três castas consideradas. A nuvem sugere também que poderá haver dispersões maiores das observações, em torno da recta de fundo, para as folhas de maior dimensão.

(b) Eis os comandos R necessários, e os resultados numéricos correspondentes:

```
> videirasN.lm <- lm(NLdir ~ NP, data=videiras)
> summary(videirasN.lm)
[...]
```

Coefficients:

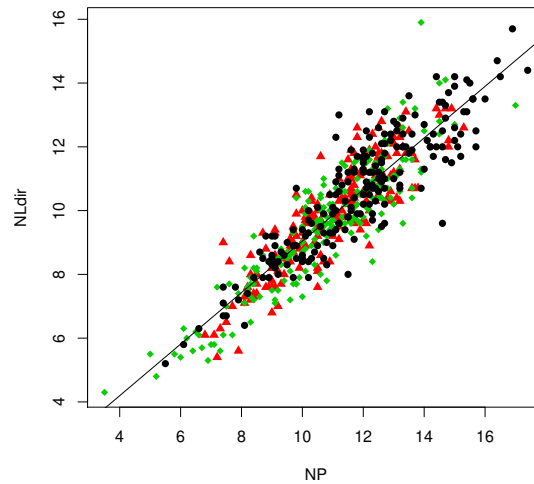
|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 0.96218  | 0.18309    | 5.255   | 2.06e-07 *** |
| NP          | 0.80841  | 0.01607    | 50.314  | < 2e-16 ***  |

```

Residual standard error: 0.8339 on 598 degrees of freedom
Multiple R-squared: 0.8089, Adjusted R-squared: 0.8086
F-statistic: 2532 on 1 and 598 DF, p-value: < 2.2e-16
```

Assim, a recta de regressão  $y = 0.96218 + 0.80841x$  explica cerca de 81% da variabilidade observada nas nervuras laterais direitas, para o conjunto das  $n = 600$  observações. Trata-se duma aproximação razoavelmente boa (como se pode constatar no gráfico), que explica cerca de 81% da variabilidade observada nas nervuras laterais direitas. Como seria de esperar, o modelo ajustado difere significativamente do modelo nulo, tendo a estatística calculada no teste  $F$  de ajustamento global um valor  $F_{calc} = 2532$ , cuja significância ( $p$ -value) correspondente é inferior à precisão de máquina, logo indistinguível de zero.

```
> abline(videirasN.lm)
```



- (c) Eis os comandos R necessários, e os resultados numéricos correspondentes ao modelo ANCOVA pedido:

```
> videirasNCasta.lm <- lm(NLdir ~ NP*Casta, data=videiras)
> summary(videirasNCasta.lm)
[...]
```

|                      | Estimate | Std. Error | t value | Pr(> t )     |
|----------------------|----------|------------|---------|--------------|
| (Intercept)          | 1.39812  | 0.32102    | 4.355   | 1.57e-05 *** |
| NP                   | 0.77780  | 0.02654    | 29.305  | < 2e-16 ***  |
| CastaFernaõ Pires    | -0.43069 | 0.48897    | -0.881  | 0.379        |
| CastaVital           | -0.66120 | 0.43788    | -1.510  | 0.132        |
| NP:CastaFernaõ Pires | 0.03395  | 0.04253    | 0.798   | 0.425        |
| NP:CastaVital        | 0.04100  | 0.03798    | 1.079   | 0.281        |

```

Residual standard error: 0.8316 on 594 degrees of freedom
Multiple R-squared: 0.8112, Adjusted R-squared: 0.8096
F-statistic: 510.5 on 5 and 594 DF, p-value: < 2.2e-16
```

A recta para a casta Água Santa (a casta correspondente ao primeiro nível do factor, o nível de referência, logo não explicitada na listagem de resultados) tem equação  $y = 1.39812 + 0.77780x$ . Para obter a equação correspondente à casta Fernão Pires, será necessário acrescentar à ordenada na origem o acréscimo estimado  $\hat{\alpha}_{0,2} = -0.43069$  e ao declive, o respectivo acréscimo estimado,  $\hat{\alpha}_{1,2} = 0.03395$ . De forma análoga, obtém-se a recta ajustada para a casta Vital. Eis as equações das três rectas ajustadas:

$$\begin{aligned} \text{Casta Água Santa} & \quad y = 1.39812 + 0.77780x \\ \text{Casta Fernão Pires} & \quad y = (1.39812 - 0.43069) + (0.77780 + 0.03395)x = 0.96743 + 0.81175x \\ \text{Casta Vital} & \quad y = (1.39812 - 0.66120) + (0.77780 + 0.04100)x = 0.73692 + 0.81880x \end{aligned}$$

Para traçar as rectas de cada casta na nuvem de pontos já criada, podem usar-se os seguintes comandos:

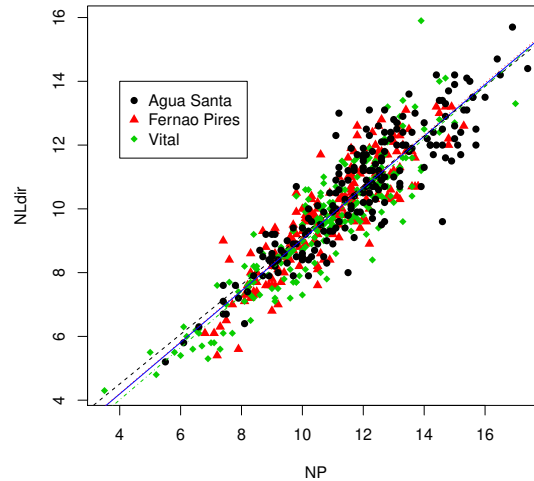
```
> coefVidCasta <- coef(videirasNCasta.lm)
> coefVidCasta
(Intercept) NP CastaFernaõ Pires CastaVital NP:CastaFernaõ Pires NP:CastaVital
1.39811600 0.77779606 -0.43068514 -0.66119902 0.03394865 0.04100268
```

```

> abline(coefVidCasta[c(1,2)], col=1, lty=2) <-- recta casta Água Santa
> abline(coefVidCasta[c(1,2)]+coefVidCasta[c(3,5)],col=2,lty=3) <-- recta casta Fernão Pires
> abline(coefVidCasta[c(1,2)]+coefVidCasta[c(4,6)],col=3,lty=4) <-- recta casta Vital

```

Apesar das equações diferentes, as quatro rectas são difíceis de distinguir no gráfico.



- (d) A equação do modelo ANCOVA ajustado pode escrever-se da seguinte forma, utilizando a notação vectorial:

$$\vec{y} = \beta_0 + \beta_1 \vec{x} + \alpha_{0:2} \vec{\mathcal{I}}_2 + \alpha_{0:3} \vec{\mathcal{I}}_3 + \alpha_{1:2} \vec{\mathcal{I}}_2 \star \vec{x} + \alpha_{1:3} \vec{\mathcal{I}}_3 \star \vec{x} + \vec{\epsilon},$$

sendo  $\vec{\mathcal{I}}_i$  a variável indicatriz das observações da casta  $i = 2, 3$  (Fernão Pires e Vital, respectivamente) e  $\alpha_{j:i}$  o acréscimo no parâmetro  $\beta_j$  (em relação à casta de referência, a Água Santa), resultante de estarmos na casta  $i = 2, 3$ . O símbolo  $\star$  indica um produto elemento a elemento entre dois vectores de igual dimensão. O modelo linear ajustado acima pode agora ser visto como um submodelo deste modelo ANCOVA, associado à hipótese  $\alpha_{0:2} = \alpha_{0:3} = \alpha_{1:2} = \alpha_{1:3} = 0$ . Vamos efectuar um teste  $F$  parcial para testar a equivalência de modelo e submodelo.

**Hipóteses:**  $H_0 : \alpha_{j:i} = 0, \forall j = 0, 1; i = 2, 3$  vs.  $H_1 : \exists j = 0, 1; i = 2, 3$  tal que  $\alpha_{j:i} \neq 0$ .

**Estatística do Teste:** (na forma mais adequada à informação disponível)

$$F = \frac{R_c^2 - R_s^2}{1 - R_c^2} \cdot \frac{n - (p+1)}{p - k} \cap F_{(p-k, n-(p+1))}, \text{ sob } H_0.$$

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica:** (unilateral direita) Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(4,594)} \approx 2.39$ .

**Conclusão:** Temos  $F_{calc} = \frac{0.8112 - 0.8089}{1 - 0.8112} \cdot \frac{594}{4} = 1.809$ . Logo, não rejeitamos  $H_0$ , isto é, não se pode dizer que o modelo ANCOVA se ajuste de forma significativamente diferente do modelo RLS com uma única recta para as três castas. Assim, não se justifica abandonar o modelo RLS, que é mais parcimonioso e tem um ajustamento considerado adequado.

Este teste  $F$  parcial, comparando o modelo ANCOVA ajustado na alínea anterior com o submodelo ajustado na alínea 1b (recta única para a totalidade das observações) obtém-se no R com o comando `anova`:

```
> anova(videirasN.lm, videirasNCasta.lm)
Analysis of Variance Table
Model 1: NLdir ~ NP
Model 2: NLdir ~ NP * Casta
 Res.Df RSS Df Sum of Sq F Pr(>F)
1 598 415.80
2 594 410.81 4 4.9948 1.8055 0.1262
```

**NOTA:** A pequena discrepância no valor calculado da estatística de teste resulta de, na nossa resolução anterior, terem sido usados valores de  $R^2$  arredondados a 4 casas decimais.

(e) Eis os três ajustamentos “mono-casta” pedidos.

- i. Tendo em atenção que as  $n_1 = 200$  observações da casta Água Santa estão nas linhas 401 a 600 da *data frame*, tem-se:

```
> summary(lm(NLdir ~ NP, data=videiras[401:600,]))
[...]
```

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 1.39812  | 0.33349    | 4.192   | 4.16e-05 *** |
| NP          | 0.77780  | 0.02757    | 28.210  | < 2e-16 ***  |

---

Residual standard error: 0.8639 on 198 degrees of freedom  
 Multiple R-squared: 0.8008, Adjusted R-squared: 0.7998  
 F-statistic: 795.8 on 1 and 198 DF, p-value: < 2.2e-16

A recta de regressão obtida ( $y = 1.39812 + 0.77780x$ ) é a mesma que no modelo completo (modelo ANCOVA) considerado acima. O valor do coeficiente de determinação ( $R^2 = 0.8008$ ) é muito próximo do valor obtido com a recta única para a totalidade das  $n = 600$  observações, facto que não era possível prever a partir dos ajustamentos anteriores.

- ii. As  $n_2 = 200$  observações da casta Fernão Pires estão nas 200 primeiras linhas do objecto *videiras*. Assim,

```
> summary(lm(NLdir ~ NP, data=videiras[1:200,]))
[...]
```

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )    |
|-------------|----------|------------|---------|-------------|
| (Intercept) | 0.96743  | 0.34914    | 2.771   | 0.00612 **  |
| NP          | 0.81174  | 0.03146    | 25.801  | < 2e-16 *** |

---

Residual standard error: 0.7872 on 198 degrees of freedom  
 Multiple R-squared: 0.7708, Adjusted R-squared: 0.7696  
 F-statistic: 665.7 on 1 and 198 DF, p-value: < 2.2e-16

Também neste caso, e como teria de ser, a recta obtida ( $y = 0.96743 + 0.81174x$ ) é, a menos de erros de arredondamento, a recta obtida ao ajustar o modelo ANCOVA. Também neste caso, o coeficiente de determinação  $R^2 = 0.7708$  é próximo do valor obtido para a recta única, embora neste caso não tenha necessariamente de ser assim.

- iii. Para as restantes observações, relativas à casta Vital, tem-se:

```
> summary(lm(NLdir ~ NP, data=videiras[201:400,]))
[...]
```

Coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )    |
|-------------|----------|------------|---------|-------------|
| (Intercept) | 0.73692  | 0.30147    | 2.444   | 0.0154 *    |
| NP          | 0.81880  | 0.02751    | 29.769  | < 2e-16 *** |



---

Residual standard error: 0.8418 on 198 degrees of freedom  
Multiple R-squared: 0.8174, Adjusted R-squared: 0.8164  
F-statistic: 886.2 on 1 and 198 DF, p-value: < 2.2e-16

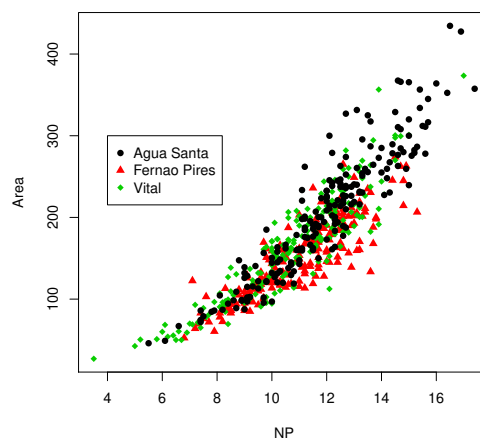
Confirma-se a recta de regressão  $y = 0.73692 + 0.8188x$ , e mais uma vez o valor  $R^2 = 0.8174$  é próximo do obtido com uma única recta de regressão para as três castas, o que é, como para as outras castas, uma particularidade deste exemplo, associada ao facto de as três nuvens de pontos serem de configuração semelhante.

- (f) O único modelo que não é de RLS é o modelo completo de ANCOVA, e será o único cuja matriz do modelo é aqui considerada. A fim de poupar no espaço, apenas se mostram as linhas correspondentes às três primeiras observações de cada casta. Recorde-se que à casta de referência (que, uma vez que o R ordena os níveis do factor por ordem alfabética, é a Água Santa) correspondem as últimas 200 linhas da matriz. As restantes castas estão indicadas nos nomes de coluna.

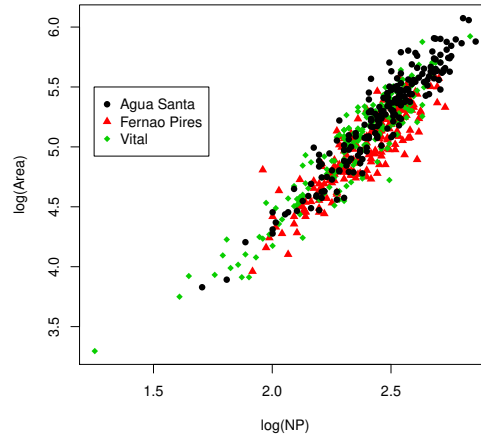
```
> model.matrix(videirasNCasta.lm)
 (Intercept) NP CastaFerna Pires CastaVital NP:CastaFerna Pires NP:CastaVital
1 1 13.8 1 0 13.8 0.0
2 1 9.1 1 0 9.1 0.0
3 1 14.5 1 0 14.5 0.0
[...]
201 1 11.7 0 1 0.0 11.7
202 1 10.6 0 1 0.0 10.6
203 1 11.0 0 1 0.0 11.0
[...]
401 1 15.7 0 0 0.0 0.0
402 1 11.7 0 0 0.0 0.0
403 1 10.2 0 0 0.0 0.0
[...]
```

2. Neste exercício a variável resposta é *Area* e a variável preditora é *NP*.

- (a) O gráfico (obtido de forma análoga ao que foi visto no Exercício 1a) torna evidente a existência duma curvatura na relação entre área foliar e comprimento da nervura principal. esta curvatura não é de estranhar, uma vez que a área é uma característica bi-dimensional, enquanto que o comprimento é unidimensional, sugerindo que a área seja aproximadamente proporcional ao quadrado do comprimento da nervura.



- (b) Com a dupla logaritmização pedida no enunciado obtém-se uma relação mais próxima da linearidade. Assim, a logaritmização de área foliar e de comprimento da nervura principal é uma boa transformação linearizante.



- (c) O modelo pedido tem o seguinte ajustamento.

```
> vid.Anc2.lm <- lm(log(Area) ~ log(NP), data=videiras)
> summary(vid.Anc2.lm)
[...]
```

|             | Estimate | Std. Error | t value | Pr(> t )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 0.57869  | 0.07703    | 7.513   | 2.12e-13 *** |
| log(NP)     | 1.87642  | 0.03203    | 58.579  | < 2e-16 ***  |

---  
Residual standard error: 0.1597 on 598 degrees of freedom  
Multiple R-squared: 0.8516, Adjusted R-squared: 0.8513  
F-statistic: 3431 on 1 and 598 DF, p-value: < 2.2e-16

A recta ajustada, às variáveis logaritimizadas é  $\ln(\text{Area}) = 0.57869 + 1.87642 \ln(\text{NP})$ . Em termos das variáveis originais (não logaritimizadas), esta relação corresponde a uma relação potência  $\text{Area} = e^{0.57869} \text{NP}^{1.87642}$  (ver acetatos das aulas relativos às transformações linearizantes).

- (d) O modelo ANCOVA agora pedido tem o seguinte ajustamento:

```
> vid.Anc2d <- lm(log(Area) ~ log(NP)*Casta, data=videiras)
> summary(vid.Anc2d)
[...]
```

|                           | Estimate | Std. Error | t value | Pr(> t )     |
|---------------------------|----------|------------|---------|--------------|
| (Intercept)               | 0.33820  | 0.13050    | 2.592   | 0.009791 **  |
| log(NP)                   | 1.99648  | 0.05294    | 37.711  | < 2e-16 ***  |
| CastaFernao Pires         | 0.62328  | 0.19914    | 3.130   | 0.001834 **  |
| CastaVital                | 0.39524  | 0.17007    | 2.324   | 0.020463 *   |
| log(NP):CastaFernao Pires | -0.31298 | 0.08232    | -3.802  | 0.000158 *** |
| log(NP):CastaVital        | -0.17654 | 0.07025    | -2.513  | 0.012232 *   |

---  
Residual standard error: 0.1482 on 594 degrees of freedom

Multiple R-squared: 0.8731, Adjusted R-squared: 0.872  
 F-statistic: 817.4 on 5 and 594 DF, p-value: < 2.2e-16

O valor do coeficiente de determinação deste modelo ( $R^2 = 0.8731$ ) é comparável com o do modelo de regressão linear simples ajustado na alínea anterior ( $R^2 = 0.8516$ ), uma vez que em ambos os casos a escala da variável resposta é a de log-áreas. O coeficiente de determinação aumentou com o modelo ANCOVA (como tem de ser, uma vez que o modelo de uma única recta de regressão é um submodelo do modelo ANCOVA), mas o aumento não é muito acentuado (pouco mais de 2%), pelo que é legítima a dúvida se o aumento obtido com o modelo ANCOVA compensa a maior complexidade do modelo.

- (e) Tendo em conta a natureza destes parâmetros estimados, resultam as seguintes relações para cada casta:

$$\begin{array}{lll} \text{Água Santa} & \ln(\text{Area}) = 0.33820 + 1.99648 \ln(NP) & \Leftrightarrow \text{Area} = e^{0.33820} NP^{1.99648} \\ \text{Fernão Pires} & \ln(\text{Area}) = 0.96148 + 1.6835 \ln(NP) & \Leftrightarrow \text{Area} = e^{0.96148} NP^{1.6835} \\ \text{Vital} & \ln(\text{Area}) = 0.73344 + 1.81994 \ln(NP) & \Leftrightarrow \text{Area} = e^{0.73344} NP^{1.81994} \end{array}$$

Em todos os casos, a área foliar é modelada como proporcional a uma potência do comprimento da nervura principal, potência essa que varia entre 1.68 e 2. Uma relação  $\text{Area} = NP^2$  corresponderia a folhas de forma quadrada, com lado igual a  $NP$ . A forma irregular da folha justifica as potências menores que 2 e as constantes de proporcionalidade, que oscilam entre 1.40 (no caso da casta Água Santa) e 2.62 (casta Fernão Pires).

- (f) Uma vez que os modelos das alíneas (c) e (d) são modelos encaixados, é possível usar um teste  $F$  parcial para estudar se o respectivo ajustamento é significativamente diferente. A equação do modelo ANCOVA é da forma

$$\vec{y} = \beta_0 + \beta_1 \vec{x} + \alpha_{0:2} \vec{\mathcal{I}}_2 + \alpha_{0:3} \vec{\mathcal{I}}_3 + \alpha_{1:2} \vec{\mathcal{I}}_2 \star \vec{x} + \alpha_{1:3} \vec{\mathcal{I}}_3 \star \vec{x} + \vec{\epsilon},$$

sendo  $\vec{\mathcal{I}}_i$  a variável indicatriz das observações da casta  $i = 2, 3$  (Fernão Pires e Vital, respectivamente) e  $\alpha_{j:i}$  o acréscimo no parâmetro  $\beta_j$  (em relação à casta de referência, a Água Santa), resultante de estarmos na casta  $i = 2, 3$ . O símbolo  $\star$  indica um produto elemento a elemento entre dois vectores de igual dimensão. O modelo linear ajustado acima pode agora ser visto como um submodelo deste modelo ANCOVA, associado à hipótese  $\alpha_{0:2} = \alpha_{0:3} = \alpha_{1:2} = \alpha_{1:3} = 0$ .

**Hipóteses:**  $H_0 : \alpha_{j:i} = 0, \forall j = 0, 1; i = 2, 3$  vs.  $H_1 : \exists j = 0, 1; i = 2, 3$  tal que  $\alpha_{j:i} \neq 0$ .

**Estatística do Teste:** (na forma mais adequada à informação disponível)

$$F = \frac{R_c^2 - R_s^2}{1 - R_c^2} \cdot \frac{n - (p+1)}{p - k} \cap F_{(p-k, n-(p+1))}, \text{ sob } H_0.$$

**Nível de significância:**  $\alpha = 0.05$ .

**Região Crítica:** (unilateral direita) Rejeitar  $H_0$  se  $F_{calc} > f_{0.05(4,594)} \approx 2.39$ .

**Conclusão:** Temos  $F_{calc} = \frac{0.8731 - 0.8516}{1 - 0.8731} \cdot \frac{594}{4} = 25.15957$ . Logo, neste caso rejeita-se claramente  $H_0$ , isto é, conclui-se que o ajustamento do modelo ANCOVA é significativamente diferente do ajustamento do modelo RLS com uma única recta para as três castas. Assim, do ponto de vista estatístico justifica-se a utilização do modelo ANCOVA, com rectas/curvas diferentes para cada casta.

O recurso ao comando `anova` do R confirma o valor calculado da estatística (arredondamentos aparte) e o valor quase nulo do  $p$ -value correspondente.

```
> anova(vid.Anc2.lm, vid.Anc2d)
Analysis of Variance Table
```

```

Model 1: log(Area) ~ log(NP)
Model 2: log(Area) ~ log(NP) * Casta
Res. Df RSS Df Sum of Sq F Pr(>F)
 1 598 15.248
 2 594 13.037 4 2.2102 25.174 < 2.2e-16 ***

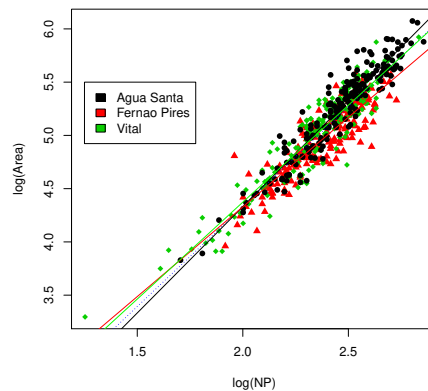
```

- (g) O gráfico pedido é indicado em baixo, sendo a recta única para a totalidade das  $n = 600$  observações indicada a tracejado. O gráfico foi construído com os seguintes comandos do R:

```

> plot(log(Area)~log(NP), col=as.numeric(Casta), pch=as.numeric(Casta)+15, data=videiras)
> abline(vid.Anc2.lm, col="blue", lty="dotted")
> abline(0.33820, 1.99648, col="black")
> abline(0.33820+0.62328, 1.99648-0.31298, col="red")
> abline(0.33820+0.39524, 1.99648-0.17654, col="green")
> legend(1.25,5.5, levels(videiras$Casta), fill=1:3)

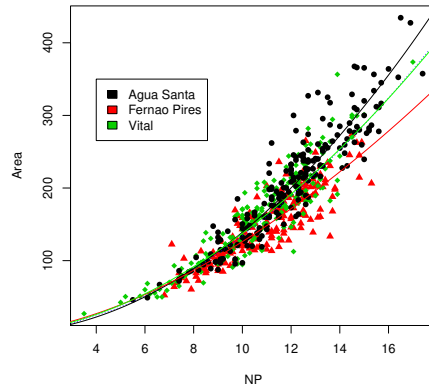
```



Confirma-se o maior declive da recta associada à casta Água Santa, e o menor associado à casta Fernão Pires. Em comparação com a relação análoga estudada no Exercício 1, é visível uma maior distinção das três rectas ajustadas, que foi reflectida no facto de o teste  $F$  parcial ter considerado que o modelo ANCOVA e o modelo de regressão linear simples para as três castas em conjunto serem significativamente diferentes.

**NOTA:** Convém acrescentar que a significância do teste  $F$  parcial resulta também do número bastante elevado de observações usado para ajustar estes modelos ( $n = 600$ ). Quanto mais informação estiver disponível na amostra, mais facilmente as diferenças são consideradas significativas.

- (h) O gráfico para as variáveis não logaritmizadas é o seguinte.



Foi produzido com os comandos:

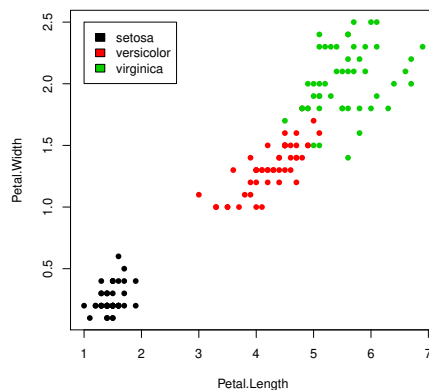
```
> plot(Area ~ NP, col=as.numeric(Casta), pch=as.numeric(Casta)+15, data=videiras)
> curve(exp(0.5787)*x^(1.8764), from=0, to=18, col="blue", lty="dotted", add=TRUE)
> curve(exp(0.3382)*x^(1.9965), from=0, to=18, add=TRUE)
> curve(exp(0.3382+0.6233)*x^(1.9965-0.3130), from=0, to=18, col="red", add=TRUE)
> curve(exp(0.3382+0.3952)*x^(1.9965-0.1765), from=0, to=18, col="green", add=TRUE)
> legend(4,350, levels(videiras$Casta), fill=1:3)
```

Nas escalas originais (não logaritmizadas) as diferenças entre as castas Água Santa e Fernão Pires é mais visível. A casta Vital tem um comportamento muito próximo do comportamento conjunto das três castas, sendo a sua curva ajustada quase indistinguível da curva única para as três castas (representada a pontado).

3. Neste exercício, consideram-se as  $n = 150$  observações sobre lírios, com variável resposta dada pela largura das pétalas (variável `Petal.Width`) e preditor numérico comprimento das pétalas (`Petal.Length`). Será considerado também o factor espécie (`Species`), havendo  $n_i = 50$  observações de cada espécie.

- (a) O gráfico pedido é obtido com os comandos seguintes. A nuvem é prometedora para uma relação linear global.

```
> plot(Petal.Width ~ Petal.Length, col=Species, data=iris, pch=16)
> legend(1,2.5, legend=levels(iris$Species), fill=1:3)
```



(b) Tem-se:

```
> iris.lm <- lm(Petal.Width ~ Petal.Length, data=iris)
> summary(iris.lm)
[...]
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.363076 0.039762 -9.131 4.7e-16 ***
Petal.Length 0.415755 0.009582 43.387 < 2e-16 ***

Residual standard error: 0.2065 on 148 degrees of freedom
Multiple R-squared: 0.9271, Adjusted R-squared: 0.9266
F-statistic: 1882 on 1 and 148 DF, p-value: < 2.2e-16
```

A recta  $y = -0.363076 + 0.415755x$  explica quase 93% da variabilidade observada nas larguras das pétalas, para o conjunto das três espécies de lírios.

(c) O modelo completo, cruzando o preditor numérico `Petal.Length` com o factor `Species` é:

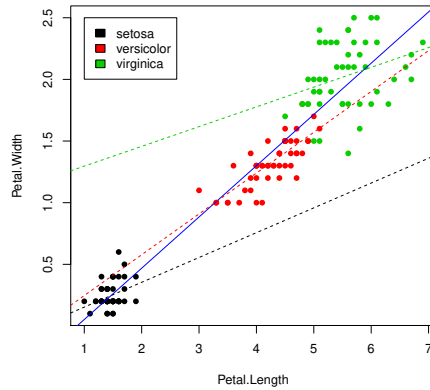
```
> irisSpecies.lm <- lm(Petal.Width ~ Petal.Length*Species, data=iris)
> summary(irisSpecies.lm)
[...]
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.04822 0.21472 -0.225 0.822627
Petal.Length 0.20125 0.14586 1.380 0.169813
Speciesversicolor -0.03607 0.31538 -0.114 0.909109
Speciesvirginica 1.18425 0.33417 3.544 0.000532 ***
Petal.Length:Speciesversicolor 0.12981 0.15550 0.835 0.405230
Petal.Length:Speciesvirginica -0.04095 0.15291 -0.268 0.789244

Residual standard error: 0.1773 on 144 degrees of freedom
Multiple R-squared: 0.9477, Adjusted R-squared: 0.9459
F-statistic: 521.9 on 5 and 144 DF, p-value: < 2.2e-16
```

- i. As três rectas de regressão, para cada espécie individual são:  $y = -0.04822 + 0.20125x$  para a espécie *setosa*,  $y = -0.08429 + 0.33106x$  para a espécie *versicolor*, e  $y = 1.1360 + 0.1603x$  para a espécie *virginica*. O valor do coeficiente de determinação do modelo ANCOVA,  $R^2 = 0.9477$  é naturalmente maior do que o  $R^2$  do submodelo constituído por uma única recta de regressão. Mas o seu valor não é de interpretação imediata, como se viu nas aulas e como se verá nas alíneas seguintes. Para traçar estas três rectas por espécie individual em cima da nuvem de pontos já anteriormente obtida, podem dar-se os seguintes comandos:

```
> coefIrisSpecies <- coef(irisSpecies.lm)
> abline(coefIrisSpecies[c(1,2)], col=1, lty=2)
> abline(coefIrisSpecies[c(1,2)]+coefIrisSpecies[c(3,5)], col=2, lty=2)
> abline(coefIrisSpecies[c(1,2)]+coefIrisSpecies[c(4,6)], col=3, lty=2)
```

Os resultados obtidos, juntamente com a recta única obtida para a totalidade das  $n = 150$  observações (a azul, em traço contínuo), são indicados no gráfico seguinte.



Como se pode constatar, a situação é bem mais confusa do que no exercício 1, com duas das rectas (das espécies *setosa* e *virginica*) com declives bastante diferentes em relação aos da recta global e da recta da espécie *versicolor*. No entanto, as rectas das espécies *setosa* e *virginica* parecem ser aproximadamente paralelas, sendo os declives ajustados (0.20125 e 0.1603) próximos. No modelo completo discutido nas aulas, o declive da recta para a espécie de referência (*setosa*) é o parâmetro  $\beta_1$ . O declive da recta para a espécie *virginica* é a soma de  $\beta_1$  com o acréscimo específico do declive da espécie *virginica*, ou seja, com o acréscimo  $\alpha_{1:3}$ . A hipótese de que essas duas rectas sejam paralelas corresponde à hipótese de  $H_0 : \alpha_{1:3} = 0$ . Esta hipótese corresponde a um teste a um parâmetro individual num modelo linear (ou seja, corresponde aos testes  $t$  usados na regressão linear para aferir possíveis valores de cada  $\beta_j$ ). A informação necessária para efectuar esse teste está disponível na listagem de resultados obtida acima para o modelo `irisSpecies.lm`. Em particular, a estimativa desse acréscimo é  $-0.04095$ , com um erro padrão associado de  $\hat{\sigma}_{\hat{\alpha}_{1:3}} = 0.15291$ . Tendo em conta a hipótese nula referida, a estatística  $t$  do teste também é dada na listagem e tem valor  $T_{calc} = -0.268$ , a que corresponde um valor de prova  $p = 0.789244$ . Sendo assim, está-se muito longe de rejeitar a hipótese nula  $H_0 : \alpha_{1:3} = 0$ , para qualquer nível de significância usual. Assim, não se rejeita que essas duas rectas de espécie são paralelas.

- (d) Os três modelos individuais de espécie, ajustados apenas usando as  $n_i = 50$  observações de cada espécie têm os coeficientes de determinação indicados de seguida:

```
> irisSetosa.lm <- lm(Petal.Width ~ Petal.Length, data=iris[1:50,])
> irisVersi.lm <- lm(Petal.Width ~ Petal.Length, data=iris[51:100,])
> irisVirgi.lm <- lm(Petal.Width ~ Petal.Length, data=iris[101:150,])
> summary(irisSetosa.lm)$r.sq
[1] 0.1099785
> summary(irisVersi.lm)$r.sq
[1] 0.6188467
> summary(irisVirgi.lm)$r.sq
[1] 0.1037537
```

Assim, em todos os casos, estes  $R^2$  por espécie individual são muito mais baixos que o  $R^2$  global correspondente ao modelo ANCOVA completo. Como se discutiu nas aulas, tal facto corresponde a uma situação em que uma ANOVA da variável resposta `Petal.Width` sobre um único factor `Species` tem um valor elevado da Soma de Quadrados correspondente ao ajustamento do modelo, ou seja,  $SQF$  elevado. Por outras palavras, o valor elevado de

$R^2 = 0.9477$  no modelo ANCOVA resulta do facto de ao factor espécie corresponderem larguras médias das pétalas bastante diferentes, e não tanto ao valor preditivo do preditor numérico `Petal.Length`. A tradução prática desse facto é visível na nuvem de pontos original, se repararmos que a forte relação linear global tem sobretudo a que ver com a separação entre os três grupos de observações correspondentes a cada espécie, e não tanto com relações lineares fortes entre as duas medições das pétalas no seio de cada espécie. Por outras palavras, a relação linear tão prometedora que parece existir entre largura e comprimento das pétalas, na nuvem da totalidade das  $n = 150$  observações, é em certo sentido uma ilusão resultante de se ter considerado em conjunto as três espécies.

- (e) Nas aulas foi vista a fórmula que relaciona o valor de  $R^2$  global do modelo ANCOVA com os  $R^2$  e as Somas de Quadrados Totais para cada subconjunto de observações (por espécie), bem como o valor de  $SQF$  na ANOVA a um factor relacionando `Petal.Width` e o factor `Species`. A fórmula é

$$R^2 = \frac{\sum_{i=1}^s R_i^2 SQT_i + SQF}{\sum_{i=1}^s SQT_i + SQF}.$$

O valor de  $SQF$  pode obter-se da seguinte forma:

```
> summary(aov(Petal.Width ~ Species, data=iris))
 Df Sum Sq Mean Sq F value Pr(>F)
Species 2 80.41 40.21 960 <2e-16 ***
Residuals 147 6.16 0.04
```

Por outro lado, os valores de  $SQT_i$  podem ser obtidos como o numerador das variâncias dos valores observados das larguras de pétalas em cada espécie. Tem-se  $SQT_1 = 49 \times s_{y_1}^2 = 0.5442$ ;  $SQT_2 = 49 \times s_{y_2}^2 = 1.9162$  e  $SQT_3 = 49 \times s_{y_3}^2 = 3.6962$ . Logo,

$$R^2 = \frac{(0.1099785 \times 0.5442) + (0.6188467 \times 1.9162) + (0.1037537 \times 3.6962) + 80.41}{(0.5442 + 1.9162 + 3.6962) + 80.41} = 0.9477001.$$

Como se pode constatar, o valor de  $SQF$  sobrepõe-se ao das restantes parcelas, quer no numerador, quer no denominador, gerando um valor muito elevado do coeficiente de determinação global do modelo ANCOVA, que não corresponde a valores elevados de  $R^2$  em nenhuma das regressões individuais de cada espécie. Confirma-se que a interpretação dos valores de  $R^2$  em modelos ANCOVA deve ser feita com cuidado.