

INSTITUTO SUPERIOR DE AGRONOMIA
ESTATÍSTICA E DELINEAMENTO – 2016-17

9 de Janeiro de 2017

Segundo Teste

Duração: 2h30

I [8,5 valores]

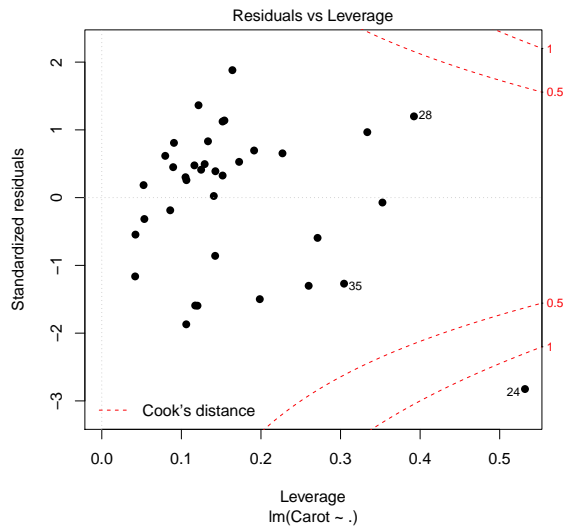
Em estudos de fisiologia da videira é importante quantificar os carotenóides. Contudo, a sua medição envolve a utilização de solventes mais caros do que os usados na medição das clorofilas e de metabólitos antioxidantes. Pretende-se estudar a modelação do teor de carotenóides (variável **Carot**, em $\mu\text{mol ml}^{-1}$), a partir das concentrações de clorofila *a* (variável **Chla**, em $\mu\text{mol ml}^{-1}$), clorofila *b* (variável **Chlb**, em $\mu\text{mol ml}^{-1}$), glutatona oxidada (variável **GSSG**, em $\mu\text{mol g}^{-1}$), glutatona reduzida (variável **GSH**, em $\mu\text{mol g}^{-1}$) e ácido ascórbico (variável **AsA**). Dispõem-se de medições destes compostos em folhas de 36 plantas da casta Trincadeira, cujas médias, desvios padrão, mínimos, máximos, e matriz de correlações são dados de seguida:

	\bar{x}	<i>s</i>	min	max		Carot	Chla	Chlb	GSSG	GSH	AsA
Carot	82.8300	55.9640	-0.158	182.915	Carot	1.000	0.989	0.883	-0.014	0.274	0.304
Chla	179.700	115.0790	2.909	385.012	Chla	0.989	1.000	0.915	-0.042	0.358	0.292
Chlb	194.800	139.2254	4.347	532.908	Chlb	0.883	0.915	1.000	0.016	0.276	0.230
GSSG	97.650	71.5406	9.034	303.103	GSSG	-0.014	-0.042	0.016	1.000	0.141	0.108
GSH	268.40	170.2865	14.975	661.843	GSH	0.274	0.358	0.276	0.141	1.000	0.228
AsA	1.2620	0.4083	0.583	2.175	AsA	0.304	0.292	0.230	0.108	0.228	1.000

1. Foi inicialmente ajustado um modelo de regressão linear múltipla utilizando a totalidade dos preditores disponíveis, obtendo-se os seguintes resultados:

```
Call: lm(formula = Carot ~ ., data = carot)
Coefficients:
            Estimate  Std. Error  t value  Pr(>|t|)
(Intercept) -4.291443    3.122114   -1.375    0.17946
Chla         0.581162    0.020504   28.344    < 2e-16
Chlb        -0.073829    0.016167   -4.567    7.90e-05
GSSG         0.042434    0.012783    3.320    0.00238
GSH         -0.037575    0.005736   -6.551    3.02e-07
AsA          2.360861    2.305387    1.024    0.31399
---
Residual standard error: ??? on 30 degrees of freedom
Multiple R-squared: 0.9925, Adjusted R-squared: 0.9913
F-statistic: 798.8 on 5 and 30 DF, p-value: < 2.2e-16
```

- (a) Calcule, justificando, uma estimativa da variância σ^2 dos erros aleatórios do modelo.
- (b) Um algoritmo de exclusão sequencial, baseado no Critério de Informação de Akaike, seleccionou um submodelo final com quatro variáveis predictoras e com $AIC = 123.64$. Diga, justificando, qual o preditor excluído. Qual o valor do Quadrado Médio Residual no submodelo escolhido? Comente.
- (c) Descreva e comente o seguinte gráfico. Tenha também em conta os valores da videira 24, indicados ao lado.



```
> t(carot[24,])
      24
Carot 103.829
Ch1a  249.572
Ch1b  101.334
GSSG   9.034
GSH   661.843
AsA    2.175
```

2. Na tentativa de obter um submodelo mais parcimonioso, foi ajustado um modelo só com dois preditores: Ch1a e GSH.

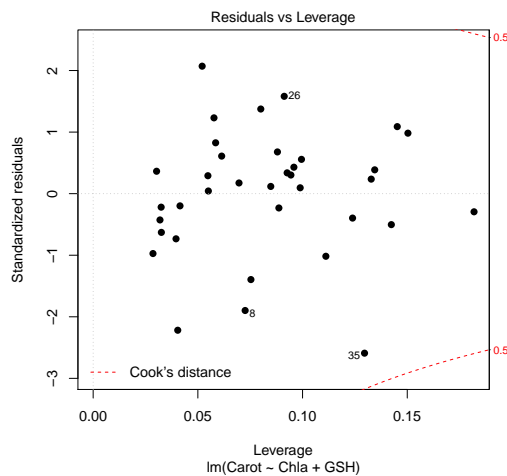
- Com base na informação disponível até aqui, indique o mais pequeno intervalo possível que garantidamente contém o valor do coeficiente de determinação desse submodelo. Comente.
- O ajustamento do referido submodelo produziu os seguintes resultados.

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.590835	2.562321	0.621	0.538959
Ch1a	0.496737	0.011120	44.672	< 2e-16
GSH	-0.029991	0.007515	-3.991	0.000345

 Residual standard error: 7.069 on 33 degrees of freedom
 Multiple R-squared: 0.985, Adjusted R-squared: 0.984
 F-statistic: 1080 on 2 and 33 DF, p-value: < 2.2e-16

- Teste formalmente ($\alpha = 0.05$) se este submodelo difere significativamente do modelo completo. Comente, tendo em conta os coeficientes de determinação de cada modelo.
- Com base na informação disponível, incluindo o gráfico indicado nesta alínea, qual dos dois modelos acima indicados escolheria? Justifique.



II [6,5 valores]

No âmbito dum estudo sobre uma variedade de arroz, pretende-se avaliar o teor final de zinco no grão (variável Zn, em mg kg⁻¹ de matéria seca), para cada um de três diferentes níveis de adubação foliar (0, 300 e 600 em mg kg⁻¹ de solo). A experiência foi realizada com base em grãos de dois tipos (polido e integral). Para cada combinação de tipo de grão e nível de adubação foliar existem 8 observações. As médias de cada um desses grupos de 8 observações são dados na tabela seguinte:

	0	300	600
grão polido	3.464	3.591	5.131
grão integral	9.934	10.313	10.664

1. Indique o delineamento experimental utilizado e descreva pormenorizadamente o modelo ANOVA adequado à experiência.
2. Um investigador efectuou uma ANOVA, tendo obtido a seguinte tabela-resumo:

	Df	Sum Sq	Mean Sq	F value
Adubacao	??	12.78	6.39	2.715
Grao	??	467.44	??	??
Adubacao:Grao	??	3.14	1.57	??
Residuals	??	??	2.35	

- (a) Complete a tabela-resumo da ANOVA, indicando como obtém os oito valores em falta
- (b) Que tipo de efeitos devem ser considerados significativos ao nível $\alpha = 0.10$? Responda de forma pormenorizada num caso, e de forma sucinta no(s) restante(s).
- (c) Que pares de combinações de tipos de grão com níveis de adubação têm médias significativamente diferentes ao abrigo da teoria de Tukey ($\alpha = 0.05$)?
- (d) Construa a tabela-resumo resultante de ajustar, aos mesmos dados, um modelo ANOVA que apenas preveja a existência de efeitos de adubação. Qual a conclusão que se teria num teste F à existência desse tipo de efeitos ($\alpha = 0.10$)? Comente.

III [5 valores]

1. Considere um modelo de regressão linear múltipla com p variáveis preditoras, ajustado com base em n observações.
 - (a) Descreva pormenorizadamente o modelo, usando a notação vectorial/matricial.
 - (b) Mostre que o vector de estimadores dos parâmetros do modelo, $\vec{\beta}$, também se pode escrever como $\vec{\beta} = \vec{\beta} + (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \vec{\epsilon}$
 - (c) Deduza a partir da expressão da alínea anterior, o vector esperado e a matriz de covariâncias do vector dos estimadores, $\vec{\beta}$, ao abrigo do modelo de regressão linear múltipla.
2. Considere os coeficientes de determinação usual (R^2) e modificado (R_{mod}^2), no contexto duma regressão linear múltipla com p variáveis preditoras, ajustada com base em n observações.
 - (a) Mostre que se verifica a relação $R_{mod}^2 = 1 - (1 - R^2) \frac{n-1}{n-(p+1)}$.

- (b) Mostre que a estatística do teste F de ajustamento global do modelo se pode escrever apenas à custa de R^2 e R_{mod}^2 , verificando-se $F_{calc} = \frac{R^2}{R^2 - R_{mod}^2}$.
- (c) Mostre que o coeficiente de determinação modificado é negativo quando $R^2 < \frac{p}{n-1}$. Comente as implicações desta condição para a estatística do teste F de ajustamento global.