

ATENÇÃO: O que se apresenta não é uma resolução, mas sim soluções, em muitos casos, sem a necessária justificação.

I [14 valores]

1. (a) A: $T_{\text{calc}} = 0.63$ (valor calculado da estatística do teste T de hipóteses $H_0 : \beta_1 = 0$ vs. $H_1 : \beta_1 \neq 0$);
 B: $R_{\text{mod}}^2 = 0.74574$ (valor do R^2 modificado);
 C: $\widehat{\text{Cov}}[\hat{\beta}_0, \hat{\beta}_0] = \hat{\sigma}_{\hat{\beta}_0}^2 = 27424.81068$ (estimativa da variância do estimador da ordenada na origem).
- (b) $R^2 = 0.7502$ indica que cerca de 75% da variância dos valores observados do teor de fenóis total é explicada por este modelo de RLM. Não é um valor muito alto, mas é significativamente diferente de zero, já que o teste F de ajustamento global ($H_0 : \mathcal{R}^2 = 0$ vs. $H_1 : \mathcal{R}^2 > 0$) apresenta um p-value $< 2.2 \times 10^{-16} \approx 0$, muito inferior a qualquer um dos usuais níveis de significância. Logo, rejeita-se a hipótese H_0 e conclui-se que o modelo difere significativamente do modelo nulo.
- (c) I. C. a 95% para β_6 : $]0.910836, 1.282464[$. Este intervalo representa os valores admissíveis do parâmetro β_6 , com 95% de confiança. Isto é, quando o teor de antocianas aumenta 1 mg/l, mantendo os restantes preditores constantes, o teor de fenóis total aumenta, em média, entre 0.910836 e 1.282464 mg/l.
- (d) A afirmação feita corresponde a $\beta_2 < -5$. Dando o ónus da prova a esta afirmação, a resposta à pergunta pode ser obtida através de um teste T com as seguintes hipóteses, $H_0 : \beta_2 \geq -5$ vs. $H_1 : \beta_2 < -5$. Depois de apresentar todos os passos deste teste, obtém-se $T_{\text{calc}} = -0.5837$ e conclui-se que não existe evidência experimental para considerar a afirmação verdadeira, ao nível de significância $\alpha = 0.05$.
- (e) A primeira variável a sair do modelo será o pH. Pela análise aos p-values dos testes T, $H_0 : \beta_i = 0$ vs. $H_1 : \beta_i \neq 0$, para $i = 1, \dots, 6$, das duas variáveis candidatas a sair do modelo (p-value $> \alpha = 0.10$), esta é a que apresenta o maior valor de p-value.
- (f)
 - i. AIC (modelo completo) = 3024.64
 - ii. O pH é a primeira variável a ser excluída pois o submodelo resultante tem menor AIC (3022.78) que o modelo inicial (completo) e que todos os restantes submodelos com menos uma variável. O peso-bago é a segunda variável a ser excluída pois o submodelo obtido com a sua exclusão apresenta um AIC (3021.27) inferior ao submodelo anterior e a todos os outros submodelos sem o pH e outra das variáveis. Mais nenhuma variável é excluída pela aplicação deste algoritmo pois os AICs dos submodelos obtidos pela saída de mais uma variável, são todos superiores ao do submodelo final do passo anterior.
 - iii. Demonstração feita na aula teórica e presente nos apontamentos do modelo linear (pág. 112).
 - iv. Teste F parcial, $H_0 : \beta_1 = \beta_3 = 0$ vs. $H_1 : \beta_1 \neq 0 \vee \beta_3 \neq 0$. Depois de apresentar todos os passos deste teste, obtém-se $F_{\text{calc}} = \frac{336}{2} \times \frac{2228748 - 2224617}{2224617} = 0.31196$ e conclui-se que modelo e submodelo não diferem significativamente, ao nível $\alpha = 0.05$, pelo que se prefere o submodelo que é mais simples.
 - v. A análise dos gráficos de resíduos apresentados deve feita considerando três aspectos:
 - breve descrição do que está representado no gráfico,
 - o que é suposto ver no gráfico se forem válidos os pressupostos do modelo linear em estudo,
 - o que está efetivamente presente no gráfico e as suas consequências para o modelo ajustado.

Os dois primeiros pontos foram apresentados nas aulas (teóricas e práticas) e encontram-se descritos nos apontamentos do modelo linear (pág. 63 a 72, 129 a 133).

Relativamente ao último ponto, o gráfico da esquerda (*Q-Q Residuals*) apresenta uma boa linearidade, indicando que o pressuposto da normalidade dos erros aleatórios deve ser válido.

O gráfico da direita (*Residuals vs Leverage*) indica que há duas observações (299 e 70) que podem ser consideradas observações atípicas (*outliers*) pois apresentam resíduos estandardizados (R_i) grandes, superiores a 4. No entanto, como os seus valores do efeito alavanca (*leverage*) são relativamente baixos e não têm uma distância de Cook elevada ($D_i < 0.5$), não são por isso consideradas observações influentes. Deste modo, ainda que possam ser sujeitas a uma análise mais detalhada para saber se podem estar associadas a algum engano, não devem gerar grande preocupação. A outra observação assinalada no gráfico (215) é a que apresenta o maior valor do efeito alavanca, cerca de 0.07, mas muito inferior ao máximo *leverage* (um). Não deve ser, por isso, considerada uma observação que tenda a "atrair" demasiado o hiperplano ajustado. Por fim, como não há pontos para além das isolinhas 0.5 da distância de Cook, não há observações que possam classificar-se de influentes, ou seja, qualquer das observações quando retirada do conjunto de dados não gera grandes alterações no ajustamento.

2. (a) $fenois = 1413.2321 - 7.2996 \text{ volumebago}$ para a região da Sardenha.
- (b) No contexto da ANCOVA, a pergunta efetuada traduz-se no teste T ao acréscimo ao declive da recta da região de referência (Saragosa) para a região Sardenha ($\alpha_{1:2}$), $H_0 : \alpha_{1:2} = 0$ vs. $H_1 : \alpha_{1:2} \neq 0$. De acordo com o $T_{\text{calc}} = -0.216$ deste teste e o correspondente p-value = 0.829, muito superior a qualquer dos habituais níveis de significância, em particular $\alpha = 0.01$, não se rejeita H_0 e podemos admitir que os declives das duas rectas não são significativamente diferentes.
- (c) Pede-se a comparação entre um modelo de covariância (rectas diferenciadas por região) e um seu submodelo (recta única). O teste F parcial, dado no comando `anova` do *output* do R, permite concluir que não se rejeita a hipótese do modelo e do submodelo não diferirem significativamente ao nível $\alpha = 0.01$, já que o correspondente p-value = 0.02931 > 0.01. Tendo os dois modelos uma qualidade idêntica (ao nível 0.01), deve preferir-se o modelo de recta única que é mais simples.
- (a) Na RLS, o coeficiente de determinação (R^2) é igual ao quadrado do coeficiente de correlação entre as variáveis preditora e resposta (r_{x^*,y^*}) e o sinal de r_{x^*,y^*} é igual ao sinal do declive da recta de regressão, logo $r_{x^*,y^*} = -0.648228$.
- (b) $fenois = 29943.27388 \times \text{volumebago}^{-0.84051}$ (relação potência).
- (c) $SQT^* = 17.278$

3. (a) Prova feita nas aulas teóricas e na página 13 dos apontamentos de apoio.

$$(b) \text{Var}(\vec{\beta}) = \begin{bmatrix} \text{Var}(\hat{\beta}_0) & \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) & \text{Var}(\hat{\beta}_1) \end{bmatrix} \text{ com as expressões dadas no formulário da UC.}$$

I [6 valores]

1. Admitindo que as observações estão ordenadas pelos respectivos níveis, $\mathbf{X}_{(9 \times 3)} = \begin{bmatrix} \vec{1_9} & \vec{I_2} & \vec{I_3} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}.$

2. (a) As unidades experimentais são as parcelas.
 (b) As pseudo-repetições são as 20 plantas de cada parcela.
 (c) Há três repetições por cada tratamento.
 (d) Em cada local, o delineamento foi totalmente casualizado.
 (e) Existem dois factores: o local (factor A com 5 níveis) e a desfolhação (factor B com 3 níveis).
 (f) Modelo ANOVA a 2 factores, factorial, com interação. Este modelo e os seus pressupostos foram descritos nas aulas teóricas e práticas, estão descritos nos slides das aulas teóricas e nos apontamentos do modelo linear (pág. 175). Deve adaptar-se a sua descrição (legenda) ao caso de estudo apresentado.
 (g) i. $A = 4$, $B = 10.44$, $C = 1452095.5$
 ii. De acordo com os resultados do teste F, há efeitos de interação local - nível de desfolhação.
 iii. Os resultados apresentados representam dois dos intervalos de confiança (IC) a 95% de Tukey:
 • IC a 95% para $\mu_{21} - \mu_{11}$: $]-3013.0314, -213.63522[$
 Como $0 \notin \text{IC}$, conclui-se, com 95% de confiança, que $\mu_{21} \neq \mu_{11}$ ou, de modo equivalente, que $\bar{y}_{21} = 2342$ e $\bar{y}_{11} = 3955$ são significativamente diferentes ao nível $\alpha = 0.05$.
 • IC a 95% para $\mu_{31} - \mu_{11}$: $]-2501.0314, 298.36478[$
 Como $0 \in \text{IC}$, conclui-se, com 95% de confiança, que $\mu_{31} = \mu_{11}$ ou, de modo equivalente, que $\bar{y}_{31} = 2854$ e $\bar{y}_{11} = 3955$ não são significativamente diferentes ao nível $\alpha = 0.05$.
3. Teste de Kruskal Wallis pois há um único factor (desfolhação).

H_0 : Os três níveis de desfolhação têm a mesma distribuição cumulativa das observações da variável resposta vs. H_1 : Os três níveis de desfolhação não têm a mesma distribuição cumulativa.

Como $H_{\text{calc}} \approx 6.49$, conclui-se que a distribuição cumulativa do número médio de sementes por planta e por parcela é diferente para os três níveis de desfolhação.