

I

Tem-se uma tabela de contingências de dimensão 4×3 , estando as linhas ($i = 1, 2, 3, 4$) associadas a variedades e as colunas ($j = 1, 2, 3$) a possíveis resultados da experiência da susceptibilidade aos fungos. O experimentador fixou o número N_i de experiências em cada variedade (sempre $N_i = 600$), pelo que apenas os totais de coluna N_j são livres.

1. O problema colocado corresponde a um teste de homogeneidade, ou seja, procura-se saber se a probabilidade de cada possível resultado da experiência (não germinados; germinados sem apressório; ou germinados com apressório) é igual, qualquer que seja a variedade, sendo esta hipótese de homogeneidade colocada como hipótese nula. Tendo sido fixados os totais das linhas, não faria sentido usar um teste de independência.

As probabilidades marginais de coluna são designadas $\pi_{.j}$ e estimadas pelas frequências relativas marginais de coluna, ou seja $\hat{\pi}_{.j} = \frac{N_{.j}}{N}$:

$$\hat{\pi}_{.1} = \frac{N_{.1}}{N} = \frac{1274}{2400} = 0.5308333 \quad ; \quad \hat{\pi}_{.2} = \frac{N_{.2}}{N} = \frac{836}{2400} = 0.34833333$$

$$\hat{\pi}_{.3} = \frac{N_{.3}}{N} = \frac{290}{2400} = 0.12083333 .$$

Ao abrigo da hipótese nula de homogeneidade, os valores esperados estimados são dados pelo produto destas probabilidades estimadas com os totais de linha, ou seja, $\hat{E}_{ij} = N_i \times \hat{\pi}_{.j} = \frac{N_i \times N_{.j}}{N}$. Como neste problema todos os totais de linha são iguais ($N_i = 600$) cada coluna tem sempre o mesmo número esperado de observações (traduzindo a homogeneidade de distribuições que é a hipótese nula ao abrigo da qual se determinam estes valores esperados estimados). Assim, tem-se, para qualquer linha $i = 1, 2, 3, 4$:

$$\hat{E}_{i1} = 600 \times 0.53083333 = 318.5 \quad ; \quad \hat{E}_{i2} = 600 \times 0.34833333 = 209.0$$

$$\hat{E}_{i3} = 600 \times 0.12083333 = 72.5 .$$

2. **Hipóteses:** Represente-se por $\pi_{j|i}$ a probabilidade do resultado $j = 1, 2, 3$, condicional a se ter a variedade $i = 1, 2, 3, 4$. As hipóteses em confronto são:

Hipótese Nula (H_0 , hipótese de homogeneidade): $\left\{ \begin{array}{l} \pi_{1|1} = \pi_{1|2} = \pi_{1|3} = \pi_{1|4} \\ \pi_{2|1} = \pi_{2|2} = \pi_{2|3} = \pi_{2|4} \\ \pi_{3|1} = \pi_{3|2} = \pi_{3|3} = \pi_{3|4} \end{array} \right.$

Hipótese Alternativa (H_1 , hipótese de heterogeneidade): pelo menos uma das desigualdades em H_0 não se verifica.

Estatística do Teste: É a estatística de Pearson, $X^2 = \sum_{i=1}^a \sum_{j=1}^b \frac{(O_{ij} - \hat{E}_{ij})^2}{\hat{E}_{ij}}$, sendo $a = 4$, $b = 3$,

O_{ij} o número de observações na célula (i, j) e \hat{E}_{ij} os valores esperados ao abrigo da hipótese de homogeneidade. A distribuição assintótica desta estatística, caso seja verdade H_0 , é $\chi_{(a-1)(b-1)}^2$ sendo $(a-1)(b-1) = 6$ graus de liberdade. É inteiramente legítimo admitir a

validade desta distribuição assintótica, uma vez que o menor dos valores esperados estimados para qualquer célula da tabela é $\hat{E}_{i3} = 72.5$ (para qualquer i), pelo que estamos bem acima do limiar de 5 (e por maioria de razão do limiar 1) referido nas condições de Cochran.

Nível de significância: $\alpha = P[\text{Erro do tipo I}] = P[\text{Rej. } H_0 \mid H_0 \text{ verdade}] = 0.05$

Região Crítica: (Unilateral direita) Rejeitar H_0 se $\chi_{\text{calc}}^2 > \chi_{\alpha[(a-1)(b-1)]}^2 = \chi_{0.05(6)}^2 = 12.592$.

Conclusões: Como $\chi_{\text{calc}}^2 = 259.7168 \gg 12.592$, tem-se uma clara rejeição da hipótese nula, ou seja, conclui-se pela rejeição da hipótese de homogeneidade.

3. Pede-se o valor, na estatística do teste, da parcela correspondente à variedade Galega ($i = 1$) e resultado "não germinado" ($j = 1$). Essa parcela é dada por:

$$\frac{(O_{11} - \hat{E}_{11})^2}{\hat{E}_{11}} = \frac{(197 - 318.5)^2}{318.5} = 46.35 .$$

Trata-se dum valor elevado, que corresponde a cerca de um quinto do valor calculado da estatística de teste, $\chi_{\text{calc}}^2 = 259.7168$, e que só por si levaria à rejeição de H_0 . Neste caso a hipótese de homogeneidade levaria a esperar um número muito superior de contagens nesta célula do que aquelas que foram efectivamente observadas. A variedade Galega é a única na qual o resultado "não germinado" não é o resultado mais frequente. Mais de dois terços dos esporos desta variedade germinaram.

II

1. Tomando recíprocos nos dois membros da equação do enunciado, obtem-se uma relação linear entre $y = \frac{1}{V}$ e $x = \frac{1}{P}$:

$$\begin{aligned} \frac{1}{V} &= \frac{1}{V_m} \cdot \frac{KP + 1}{KP} = \frac{1}{V_m} \left(1 + \frac{1}{K} \frac{1}{P} \right) \\ \Leftrightarrow \underbrace{\frac{1}{V}}_{=y} &= \underbrace{\frac{1}{V_m}}_{=\beta_0} + \underbrace{\frac{1}{V_m K}}_{=\beta_1} \cdot \underbrace{\frac{1}{P}}_{=x} \\ \Leftrightarrow y &= \beta_0 + \beta_1 x \end{aligned}$$

2. A recta de regressão ajustada entre y e x tem equação $y = 1.98780 + 4.14782 x$.

(a) No ponto anterior viu-se que, na recta de regressão entre y e x , tem-se $\beta_0 = \frac{1}{V_m}$, ou seja, $V_m = \frac{1}{\beta_0}$. Tendo em conta a ordenada na origem da recta ajustada, o valor estimado para V_m é $\hat{V}_m = \frac{1}{\hat{\beta}_0} = \frac{1}{1.98780} = 0.5030687$. De forma análoga, tem-se que o declive da recta é $\beta_1 = \frac{1}{V_m K} = \beta_0 \frac{1}{K}$. Logo, $K = \frac{\beta_0}{\beta_1}$. O correspondente valor estimado é $\hat{K} = \frac{\hat{\beta}_0}{\hat{\beta}_1} = \frac{1.98780}{4.14782} = 0.479236$. Substituindo na expressão da curva de Langmuir obtem-se a equação da curva ajustada: $V = \frac{0.503687 \times 0.4792396 \times P}{0.4792396 \times P + 1}$.

(b) Pode usar-se a recta entre y e x para obter um intervalo de predição (95%) para $y = \frac{1}{V}$ quando $P = 10$ Torr, ou seja, quando $x = \frac{1}{10} = 0.1$. A expressão geral desse intervalo de predição encontra-se no formulário, e os seus valores extremos são dados por:

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

Tem-se $b_0 + b_1 x = 1.98780 + 4.14782 \times 0.1 = 2.402672$; $n = 6$; $t_{0.025(4)} = 2.77645$; $\sqrt{QMRE} = 0.2287$. É preciso ter cuidado com os valores da média e variância de x , pois trata-se da média e variância dos *recíprocos* da pressão. A forma mais fácil de obter esses valores é utilizar as expressões que definem os parâmetros da recta: $b_1 = \frac{Cov_{xy}}{s_x^2}$ e $b_0 = \bar{y} - b_1 \bar{x}$, onde surgem esses valores, juntamente com os valores de $\bar{y} = 7.02348$ e $Cov_{xy} = 16.82589$, dados no enunciado. Da expressão para b_1 sai $s_x^2 = \frac{Cov_{xy}}{b_1} = \frac{16.82589}{4.14782} = 4.0565622$. Da expressão para b_0 sai $\bar{x} = \frac{\bar{y} - b_0}{b_1} = \frac{7.02348 - 1.98780}{4.14782} = 1.2140546$. Substituindo os valores na fórmula obtém-se o intervalo $] 1.6989755, 3.1061885 [$, que é um intervalo de predição *para o recíproco do volume*, $\frac{1}{V}$, quando $P = 10$ Torr. Tomando os recíprocos destes extremos (e necessariamente trocando a sua ordem), obtém-se um intervalo de predição (a 95%) *para o volume*, quando $P = 10$ Torr: $] 0.32194, 0.58859 [$. Esta gama de valores é coerente com os valores observados na tabela de dados, para pressões próximas de $P = 10$.

III

1. Em qualquer modelo linear, os graus de liberdade residuais são dados pela diferença entre o número total de observações e o número de parâmetros do modelo. No modelo de regressão linear múltipla em apreço, tem-se $g.l.(QMRE) = n - (p + 1) = 28$. Uma vez que existem $p = 8$ variáveis preditoras, tem-se que o ajustamento foi feito com base em $n = 28 + 9 = 37$ observações.
2. Em geral, os coeficientes β_j associados a variáveis preditoras numa regressão linear múltipla representam a variação no valor esperado da variável resposta Y , associado a um aumento duma unidade na respectiva variável preditora X_j , e mantendo constantes os valores dos restantes preditores. Assim, a estimativa $b_1 = -14.11886$ corresponderia a dizer que a transpiração diária média diminui $14.11886 \text{ mm}/\text{dia}$, caso o preditor BLUE aumente uma unidade (para valores fixos dos restantes preditores). Mas, como é sublinhado no enunciado, no nosso caso não se verificam aumentos de uma unidade nesse preditor, tendo em conta que a gama de valores desse preditor associados ao problema varia entre 0.04 e 0.12. Pode repetir-se o raciocínio que conduziu à interpretação acima referida de β_j , mas considerando aumentos de k unidades em X_j (para qualquer constante real k). Nesse caso, a diferença nos valores esperados de Y , após e antes um tal aumento (e em igualdade das restantes circunstâncias) é dado por:

$$\begin{aligned}
 E[Y | x_1, \dots, x_j + k, \dots, x_p] &= \cancel{\beta_0} + \cancel{\beta_1 x_1} + \dots + \cancel{\beta_{j-1} x_{j-1}} + \underbrace{\beta_j(x_j + k)}_{=\beta_j x_j + k \beta_j} + \cancel{\beta_{j+1} x_{j+1}} + \dots + \cancel{\beta_p x_p} \\
 - \quad E[Y | x_1, \dots, x_j, \dots, x_p] &= \cancel{\beta_0} + \cancel{\beta_1 x_1} + \dots + \cancel{\beta_{j-1} x_{j-1}} + \cancel{\beta_j x_j} + \cancel{\beta_{j+1} x_{j+1}} + \dots + \cancel{\beta_p x_p} \\
 \hline
 &= k \beta_j
 \end{aligned}$$

Tomando, por exemplo, $k = 0.01$, podemos afirmar que a estimativa da variação esperada na média da transpiração diária, quando a reflectância BLUE aumenta 0.01 unidades (mantendo outros preditores constantes) é $k b_1 = -0.1411886 \text{ mm}/\text{dia}$, ou por outras palavras, a transpiração diária média diminui $0.1411886 \text{ mm}/\text{dia}$.

3. Pede-se um teste em que a hipótese alternativa seja $\beta_8 > 0$. Tem-se o seguinte teste t :

Hipóteses: $H_0 : \beta_8 \leq 0$ vs. $H_1 : \beta_8 > 0$.

Estatística do Teste: $T = \frac{\hat{\beta}_8 - \beta_{8|H_0}}{\hat{\sigma}_{\hat{\beta}_8}} \cap t_{n-(p+1)}$, sendo H_0 verdade.

Nível de significância: $\alpha = P[\text{Erro do tipo I}] = P[\text{Rej. } H_0 \mid H_0 \text{ verdade}] = 0.05$.

Região Crítica: (Unilateral direita) Rejeitar H_0 se $T_{calc} > t_{\alpha(28)} = 1.70113$.

Conclusões: O valor de T_{calc} é dado na listagem, sendo $T_{calc} = \frac{0.04466-0}{0.02403} = 1.858 > 1.70113$.

Logo, rejeita-se H_0 e é possível concluir pela afirmação do enunciado: aumentos de temperatura estão associados a aumentos de transpiração diária (sendo as restantes condições observadas constantes). Note-se, no entanto, que não se pode afirmar que aumentos de temperatura *provoquem* aumentos de transpiração: a associação estatística não é sinónimo duma relação de causa e efeito.

4. É dado um submodelo com apenas $k=4$ preditores.

(a) Pede-se um teste F parcial para comparar o submodelo com o modelo completo original, de $p=8$ preditores. Tem-se:

Hipóteses: $H_0 : \mathcal{R}_c^2 = \mathcal{R}_s^2$ vs. $H_1 : \mathcal{R}_c^2 > \mathcal{R}_s^2$, onde \mathcal{R}_c^2 e \mathcal{R}_s^2 indicam os coeficientes de determinação populacional, respectivamente do modelo completo e do submodelo.

Estatística do Teste: $F = \frac{n-(p+1)}{p-k} \frac{\mathcal{R}_c^2 - \mathcal{R}_s^2}{1 - \mathcal{R}_c^2} \cap F_{(p-k, n-(p+1))}$, sob H_0 .

Nível de significância: $\alpha = P[\text{Erro do tipo I}] = P[\text{Rej. } H_0 \mid H_0 \text{ verdade}] = 0.05$.

Região Crítica: (Unilateral direita) Rejeitar H_0 se $F_{calc} > f_{0.05[4,28]} = 2.71$.

Conclusões: Tem-se $F_{calc} = \frac{28}{4} \frac{0.9032 - 0.8746}{1 - 0.9032} = 2.068$. Logo, não se rejeita H_0 , i.e., considera-se que a qualidade de ajustamento do modelo completo não é significativamente melhor (ao nível $\alpha=0.05$) do que o da do submodelo, que é assim um modelo mais parcimonioso. A perda de menos de 3% na variabilidade explicada da variável resposta parece ser compensada pela redução em metade no número de variáveis preditoras.

(b) As variáveis BLUE e RED formam, juntamente com o preditor GREEN, um grupo de variáveis muito fortemente correlacionadas entre si. A mais baixa correlação entre pares destas três variáveis é 0.97. Assim, é natural que a informação presente em dois destes três preditores seja bem substituída pelo terceiro destes preditores, sem grande prejuízo no valor de R^2 .

(c) Por definição, um resíduo é dado pela diferença entre um valor observado y_i e o correspondente valor ajustado através da regressão linear, \hat{y}_i . Para a sétima observação, e tendo em conta os valores dados no enunciado, tem-se:

$$\begin{aligned} e_7 &= y_7 - \hat{y}_7 = y_7 - (b_0 + b_1 \text{GREEN}_7 + b_2 \text{NIR}_7 + b_3 \text{NDVI}_7 + b_4 \text{TIR}_7) \\ &= 0.8437172 - (-3.64659 + 66.82066 \times 0.0791 - 24.96892 \times 0.1839 + \\ &\quad + 9.04922 \times 0.3399 + 0.05009 \times 17.05) \\ &= -0.1333 \end{aligned}$$

IV

1. Trata-se dum delineamento factorial a dois factores, sendo a variável resposta Y o rendimento (em kg/planta); o primeiro factor (A) o genótipo, com $a=12$ níveis (uma vez que os correspondentes graus de liberdade, indicados na tabela resumo, são $a-1=11$); e o segundo factor (B) os anos, com $b=5$ níveis (referidos no enunciado). A estrutura da tabela confirma tratar-se dum delineamento

factorial, tendo sido ajustado o modelo que prevê efeitos de interacção. Neste modelo, os graus de liberdade são dados pela diferença entre o número total de observações (n) e o número de células (situações experimentais), $ab = 12 \times 5 = 60$. Como a tabela indica que $n - ab = 420$, tem-se ao todo $n = 480$ observações. Tratando-se dum delineamento equilibrado (como é referido no enunciado), houve $n_c = \frac{n}{ab} = \frac{480}{60} = 8$ observações em cada célula. O modelo ajustado é:

- O rendimento da k -ésima parcela associada ao genótipo i , no ano j , é dado por $Y_{ijk} = \mu_{11} + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}$, $\forall i, j, k$, sendo μ_{11} o rendimento esperado do primeiro genótipo em 1994; α_i o efeito principal (acréscimo ao rendimento) associado ao genótipo i (com a restrição $\alpha_1 = 0$); β_j o efeito principal (acréscimo ao rendimento) associado ao ano j (com a restrição $\beta_1 = 0$); $(\alpha\beta)_{ij}$ o efeito de interacção (acréscimo ao rendimento) associado à combinação do genótipo i com o ano j (e com as restrições $(\alpha\beta)_{ij} = 0$ se $i = 1$ ou $j = 1$); e finalmente ϵ_{ijk} o erro aleatório da referida observação.
 - 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 ϵ_{ijk} são independentes.
2. A variância dos erros aleatórios é, pelo segundo ponto do modelo, dada por $V[\epsilon_{ijk}] = \sigma^2$. Em qualquer modelo linear, uma tal variância é estimada pelo Quadrado Médio Residual. Logo, $\hat{\sigma}^2 = QMRE = \frac{198.44}{420} = 0.4725$. As unidades de medida dos resíduos são iguais às unidades de medida da variável resposta, e tendo em conta que esses resíduos são elevados ao quadrado no cálculo de $SQRE$, tem-se que as unidades de medida desta QMRE são (kg/planta)².
3. Vai-se efectuar em pormenor o teste aos efeitos principais do Factor B (Ano), e descrever sinteticamente os testes aos efeitos principais do Factor A (genótipo) e aos efeitos de interacção.

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

Estatística do Teste: $F_B = \frac{QMB}{QMRE} \cap F_{[(b-1), n-ab]}$, sob H_0 .

Nível de significância: $\alpha = 0.05$.

Região Crítica: (Unilateral direita) Rejeitar H_0 se $F_{calc} > f_{0.05(4,420)} \approx 2.4$ (entre os valores tabelados 2.37 e 2.45).

Conclusões: Como $F_{calc} = \frac{QMB}{QMRE} = 110.180 \gg 2.45$, rejeita-se claramente H_0 , sendo possível concluir pela existência de efeitos principais de ano (ao nível $\alpha = 0.05$ e, presumivelmente para todos os níveis usuais de α). No nosso contexto, tal corresponde a afirmar que, em pelo menos alguns anos, há uma mudança no rendimento médio, quando comparado com outros anos.

No teste aos efeitos de interacção, com hipóteses $H_0 : (\alpha\beta)_{ij} = 0$, para todo o i e j , contra $H_1 : \text{existe pelo menos uma célula } (i, j) \text{ onde } (\alpha\beta)_{ij} \neq 0$, o p -value muito elevado ($p = 0.354688$) indica a não rejeição de H_0 para os habituais níveis de significância, pelo que se pode concluir pela inexistência de efeitos significativos de interacção.

Finalmente, no teste aos efeitos principais do factor genótipo, com hipóteses $H_0 : \alpha_i = 0, \forall i$ vs. $H_1 : \exists i$ tal que $\alpha_i \neq 0$, tem-se um valor de prova muito baixo ($p = 0.000904$) e para qualquer dos níveis de significância usuais rejeita-se H_0 , ou seja, conclui-se pela existência de efeitos principais de genótipo nos rendimentos médios. Assim, vale a pena escolher (ou excluir) alguns dos genótipos estudados.

4. A tabela-resumo da ANOVA correspondente ao modelo com o único factor, genótipo (ou seja, ao Modelo M_A), tem apenas 2 linhas: a correspondente à variabilidade associada ao factor e a correspondente à variabilidade residual. Por definição, a Soma de Quadrados, grau de liberdade e, por conseguinte, o Quadrado Médio associado ao factor genótipo são calculados como na tabela-resumo do modelo ANOVA a dois factores, com efeitos de interacção (o modelo M_{A*B} , onde o factor genótipo desempenhava o papel de Factor A). Logo, os correspondentes valores são iguais nas duas tabelas. Uma vez que a soma de todas as Somas de Quadrados em cada modelo ANOVA tem de ser sempre igual a $SQT = (n-1)s_y^2$, e uma vez que os valores da variável resposta Y_i com que se ajusta os dois modelos são os mesmos, tem de verificar-se $SQRE_A = SQB + SQAB + SQRE_{A*B} = 208.25 + 22.29 + 198.44 = 428.98$. De forma análoga, os graus de liberdade das duas tabelas têm de somar $n-1$, pelo que $g.l.(SQRE_A) = 4 + 44 + 420 = 468$. Tem-se então $QMRE = \frac{SQRE}{n-k} = \frac{428.98}{468} = 0.9166239$ e $F = \frac{QMF}{QMRE} = \frac{1.39}{0.9166239} = 1.51643$. A tabela-resumo vem assim:

	g.l.	SQs	QMs	F_{calc}
Factor genótipo	11	15.31	1.39	1.51643
Residual	468	428.98	0.9166239	

A principal conclusão a extrair é que ao nível de significância $\alpha = 0.05$, os efeitos de genótipo seriam agora considerados *não significativos*, uma vez que o limiar da região crítica seria agora $f_{0.05(11,468)}$ a que, pelas tabelas, corresponde um valor entre 1.91 e 1.75. Assim, os efeitos de genótipo deixariam de ser considerados significativos a esse nível de α . O facto de ignorar a existência de efeitos presentes na realidade levou ao inflacionamento da variabilidade residual, e consequentemente à diminuição do valor da estatística F , apesar de a variabilidade associada ao factor genótipo ter permanecido igual.

V

1. A matriz quadrada \mathbf{H} ser simétrica significa que $\mathbf{H}^t = \mathbf{H}$. Ora, tendo em conta as propriedades de produtos matriciais (também constantes do formulário, nomeadamente $(\mathbf{AB})^t = \mathbf{B}^t \mathbf{A}^t$; $(\mathbf{A}^t)^t = \mathbf{A}$; e $(\mathbf{A}^{-1})^t = (\mathbf{A}^t)^{-1}$), tem-se:

$$\begin{aligned} \mathbf{H}^t &= [\mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t]^t = (\mathbf{X}^t)^t [(\mathbf{X}^t \mathbf{X})^{-1}]^t \mathbf{X}^t = \mathbf{X} [(\mathbf{X}^t \mathbf{X})^t]^{-1} \mathbf{X}^t \\ &= \mathbf{X} [\mathbf{X}^t (\mathbf{X}^t)^t]^{-1} \mathbf{X}^t = \mathbf{X} (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t = \mathbf{H} , \end{aligned}$$

como se queria mostrar. Por outro lado, \mathbf{H} ser idempotente significa que $\mathbf{HH} = \mathbf{H}$. Ora,

$$\mathbf{HH} = \mathbf{X} \underbrace{(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \cdot \mathbf{X} (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t}_{=\mathbf{I}_m} = \mathbf{X} (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t = \mathbf{H} .$$

2. Tendo em conta que, no contexto duma ANOVA a um factor, a tradicional Soma de Quadrados associada ao ajustamento do modelo (que na regressão linear se designa SQR) é chamada SQF , tem-se $R^2 = \frac{SQF}{SQT}$.

- (a) A condição $R^2 = 0$ equivale a $SQF = 0$. Ora, no contexto ANOVA a um factor tem-se (ver formulário e tendo em conta que o delineamento é equilibrado):

$$SQF = \sum_{i=1}^k n_i (\bar{Y}_{i.} - \bar{Y}_{..})^2 = n_c \sum_{i=1}^k (\bar{Y}_{i.} - \bar{Y}_{..})^2 = 0 .$$

Ora, uma soma de quadrados só se pode anular se *todas* as suas parcelas se anulam o que, neste contexto, significa que $\bar{Y}_i = \bar{Y}_{..}$, para todo o i . Por outras palavras, $R^2 = 0$ se e só se todas as médias amostrais de nível forem iguais à média amostral da totalidade das observações (e portanto iguais entre si). Assim, a informação proveniente da amostra aponta de forma clara em abono da hipótese de igualdade de todas as médias populacionais de nível ($\mu_1 = \mu_2 = \dots = \mu_k$), que é a hipótese nula no teste F dum ANOVA a um único factor. Este resultado é inteiramente coerente com a não rejeição da hipótese nula do teste que resulta do facto de $R^2 = 0 \Leftrightarrow F_{calc} = 0$. Repare-se ainda que a condição $SQF = 0$ é equivalente a dizer que $SQT = SQF + SQRE = SQRE$, ou seja, toda a variabilidade de Y é residual, ou seja, interna aos níveis do factor.

- (b) A condição $R^2 = 1$ equivale a $SQF = SQT$, ou seja, $SQRE = 0$. Ora, no contexto ANOVA a um factor tem-se (ver formulário e para um delineamento equilibrado):

$$SQRE = \sum_{i=1}^k (n_i - 1) S_i^2 = (n_c - 1) \sum_{i=1}^k S_i^2 = 0 .$$

De novo, uma soma de quadrados só pode ser nula se *todas* as suas parcelas forem nulas, pelo que $SQRE = 0$ equivale a $S_i^2 = 0$, para todo o nível i , ou seja, não existe variabilidade das observações de Y no seio dum mesmo nível do factor. Neste caso tem-se também $QMRE = \frac{SQRE}{n-k} = 0$. Embora não seja possível construir a estatística do teste $F = \frac{QMF}{QMRE}$, a divisão por zero sugere um valor limite infinito, que corresponderia sempre à rejeição da hipótese nula de igualdade das médias populacionais de nível μ_i , o que é coerente com o referido facto de, neste caso, toda a variabilidade nas observações de Y corresponder à mudança entre níveis do factor.

3. Considere-se o contexto dum regressão linear simples.

- (a) Por definição, a Soma de Quadrados associada à regressão é dada por $SQR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$.
 Numa regressão linear simples, sabemos que $\hat{y}_i = b_0 + b_1 x_i$ e que (pela definição do parâmetro b_0) $\bar{y} = b_0 + b_1 \bar{x}$. Substituindo na expressão para SQR , tem-se $SQR = \sum_{i=1}^n [(b_0 + b_1 x_i) - (b_0 + b_1 \bar{x})]^2 = \sum_{i=1}^n [b_1(x_i - \bar{x})]^2 = b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 = b_1^2 (n-1) s_x^2$, como se queria mostrar.
- (b) Por definição, $R^2 = \frac{SQR}{SQT}$, logo $SQR = R^2 SQT$. Ora, em qualquer modelo linear, $SQT = (n-1) s_y^2$ e, no caso dum regressão linear simples, o coeficiente de determinação R^2 é o quadrado do coeficiente de correlação entre x e y , pelo que $SQR = r_{xy}^2 (n-1) s_y^2$, como se queria mostrar.