

Considere o arquivo `congenitos.csv`, $Y_{35.129 \times 50}$ que corresponde a dados de expressão gênica (\log_2 transformados) de 35.129 fragmentos genéticos do cromossomo de ratos congênitos. No experimento realizado, 5 réplicas de 10 ratos foram selecionadas aleatoriamente das linhagens congênicas A, B, C, D e E, sendo 2 de cada. Então, dentro de cada réplica, metade foi aleatoriamente selecionada para receber uma dieta rica em sal e a outra metade recebeu dieta controle. A Tabela 1 ilustra o esquema do experimento realizado, o qual seguiu um Fatorial 2x5 com 5 réplicas, em que, em cada uma das 50 caselas ($50=2 \times 5 \times 5$), a expressão gênica dos 35.129 fragmentos foi avaliada. Note que se trata de uma matriz com $n=50$ unidades amostrais estruturadas de acordo com o esquema fatorial 2x5 (dois fatores: Dieta (em dois níveis, Sal e Controle) e Linhagem de rato (em cinco níveis, A, B, C, D e E)) e 5 réplicas, com $p=35.129$ variáveis sendo avaliadas. Além disso, a Tabela 1 indica uma estrutura de blocagem (nas réplicas) em que, materiais genéticos de cada grupo de 10 animais (um de cada dos 10 tratamentos), foram aleatoriamente colocados juntos em *plates* para a leitura das expressões.

Tabela 1. Esquema de delineamento Fatorial 2x5 aleatorizado em blocos (réplicas)

Réplica	Sal					Controle				
	A	B	C	D	E	A	B	C	D	E
R1										
R2										
R3										
R4										
R5										

1. Escolha as Top 10 variáveis dentre as $p=35.129$. Para tanto, defina algum critério, por exemplo, seleção aleatória, ou, de forma mais interessante, você pode realizar p ANOVAs e selecionar os genes com expressão mais significativa. Note, que na ANOVA há 9 graus de liberdade para estudar o efeito de tratamento (efeito principal de Sal=1 g.l.; efeito principal de linhagem=4 g.l.; efeito de interação=4 g.l), além de 4 graus de liberdade para bloco e o restante para o termo de erro.

Também, se quiser realizar uma solução mais sofisticada ainda, e ganhar até 1 ponto a mais nesta lista, represente os 35.129 p-valores da ANOVA em um Gráfico Vulcão. O que é isso? É um gráfico de dispersão com $-\log_{10}(p\text{-valor})$ na ordenada e uma medida de tamanho de efeito (*change-fold*) na abscissa, em geral, representada pela estimativa de um efeito ($\hat{\beta}$). Pesquise sobre esse gráfico (*Vulcano plot*) e implemente.

2. Com base na seleção dos Top 10 genes, será obtida uma matriz de trabalho $Y_{50 \times 10}$ ($n=50$, $p=10$). Com base na sua matriz de trabalho represente os dados por meio de dois eixos principais obtidos de:

- Análise de Componentes Principais.
- Escalonamento Multidimensional
- Análise de Fatores (comuns e específicos)

d) Análise Discriminante: neste caso você pode adotar 2 grupos (Sal e Controle), 5 grupos (A, B, C, D e E), ou 10 grupos (do esquema fatorial 2x5). Faça sua escolha e justifique.

Em cada caso apresente as suposições envolvidas e interprete o padrão de variação mostrado nos gráficos.

3. Considere a análise discriminante linear de Fisher de sua matriz de trabalho, $Y_{50 \times 10}$, adotando $G=2$ grupos (Dieta de Sal e Controle). Além disso, considere também a análise discriminante via regressão logística (binária). Em cada caso obtenha a matriz de confusão via o método empírico e validação cruzada (*leave-one-out*). Compare os dois métodos de classificação dos grupos relativamente às suposições envolvidas e métricas de predição.

4. Com base na sua matriz de trabalho, $Y_{50 \times 10}$, obtenha a tabela de MANOVA. Há efeito significativo de algum fator sob estudo? De acordo com a MANOVA é possível obter uma decomposição da matriz de soma de quadrados e produto cruzados total (SQPC total). De maneira correspondente, obtenha uma decomposição da matriz de trabalho $Y_{50 \times 10}$. Qual é a utilidade desse tipo de decomposição na análise de dados? Como exemplo, pense nas análises realizadas no item 2, as quais usaram a matriz de trabalho mas poderiam adotar componentes dos dados.

5. A partir de uma matriz de dados normalizados $Y_{n \times p}^*$, considere a matriz de covariâncias $nS_{p \times p} = Y^* Y^{*'} = V \Lambda V'$, tal que $V_{p \times p} = (V_1, \dots, V_p)$ e $\Lambda = \text{diag}(\lambda_j)$ são matrizes de autovetores (das colunas de $Y_{n \times p}^*$) e autovalores, respectivamente, e a matriz de distâncias $D_{n \times n}$, tal que seus elementos são função dos elementos de $B_{n \times n} = Y^* Y^{*'} = U \Lambda U'$, com $U_{n \times n} = (U_1, \dots, U_n)$ matriz de autovetores (das linhas de $Y_{n \times p}^*$). Três pesquisadores realizaram análises estatísticas e chegaram à seguinte conclusão:

Pesquisador 1: $Y_{n \times p}^* \cong Y^* (V_1 \ V_2)$, $\Sigma \cong V_1 V_1' \lambda_1 + V_2 V_2' \lambda_2$

Pesquisador 2: $Y_{n \times p}^* \cong Y^* \begin{pmatrix} \frac{V_1}{\sqrt{\lambda_1}} & \frac{V_2}{\sqrt{\lambda_2}} \end{pmatrix}$

$\Sigma \cong \Phi \Phi' + \Psi$; $\Phi = (V_1 \sqrt{\lambda_1} \ V_2 \sqrt{\lambda_2})$, $\Psi = \text{diag}(1 - (\phi_1^2 + \phi_2^2))$

Pesquisador 3: $Y_{n \times p}^* \cong (U_1 \sqrt{\lambda_1} \ U_2 \sqrt{\lambda_2})$, $B \cong U_1 U_1' \lambda_1 + U_2 U_2' \lambda_2$

Qual análise estatística cada pesquisador realizou? Eles partiram do mesmo objetivo? Se desejar ganhar até 2 pontos extras nesta lista, simule dados e realize as análises dos três pesquisadores.

Boa sorte 😊