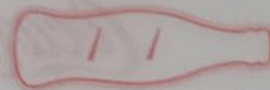


Questão 1

Coca-Cola



Jessica Caroline A. N. Temporal 7547611

Questão 1

Seja $\xi = r_1 x_1 + r_2 x_2 + \dots + r_n x_n$

$$\begin{aligned} \text{temos que } \text{Var}(\xi) &= \text{Var}(r_1 x_1 + r_2 x_2 + \dots + r_n x_n) \\ &= \text{Var}\left(\sum_{j=1}^n r_j x_j\right) \\ &= \text{Var}\left(\sum_{j=1}^n r_j x_j\right) \end{aligned}$$

Lembrando a propriedade:

$$\text{Var}(ax + by) = a^2 \text{Var}(x) + b^2 \text{Var}(y) + 2ab \text{Cov}(x, y)$$

e que como visto acima

$$\text{Var}\left(\sum_{j=1}^n r_j x_j\right) = \text{Var}(r_1 x_1 + r_2 x_2 + \dots + r_n x_n)$$

temos que:

$$\begin{aligned} \text{Var}\left(\sum_{j=1}^n r_j x_j\right) &= \text{Var}(r_1 x_1) + \text{Var}(r_2 x_2) + 2r_1 r_2 \text{Cov}(x_1, x_2) + \\ &\quad \text{Var}(r_3 x_3) + \text{Var}(r_4 x_4) + 2r_3 r_4 \text{Cov}(x_3, x_4) + \\ &\quad \vdots \\ (I) \quad &\text{Var}(r_1 x_1) + \text{Var}(r_n x_n) + 2r_1 r_n \text{Cov}(x_1, x_n) + \\ &\quad \text{Var}(r_2 x_2) + \text{Var}(r_n x_n) + 2r_2 r_n \text{Cov}(x_2, x_n) + \\ &\quad \vdots \\ &\text{Var}(r_n x_n) + \text{Var}(r_n x_n) + 2r_n r_n \text{Cov}(x_n, x_n) \end{aligned}$$

resumindo: para n dimensões temos $\frac{n!}{(n-2)! \cdot 2!}$ valores para covariância

reescrivendo de uma forma "menor" (J) temos:

$$\text{Var} \left(\sum_{j=1}^n a_j x_j \right) = \sum_{j=1}^n a_j^2 \text{Var}(x_j) + 2 \sum_{\substack{i=1, \dots, n \\ j=1, \dots, n \\ i < j}} a_i a_j \text{Cov}(x_i, x_j)$$

↳ Interpretando em forma matricial:

$$q_{jk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$$

temos um vetor $\mathbf{x} = [x_1, x_2, \dots, x_n]$

então uma matriz de k linhas e n colunas é dada por:

$$\mathbf{Q} = \frac{1}{n-1} (\mathbf{X} - \mathbf{X}_n^1)^T (\mathbf{X} - \mathbf{X}_n^1)^T$$

ou $\sum_{i=1}^n a_i x_i \text{Cov}(x_i, x_j)$ na forma matricial:

$$\begin{bmatrix} a_1 a_1 \text{Cov}(x_1, x_1) + a_1 a_2 \text{Cov}(x_1, x_2) + \dots + a_1 a_n \text{Cov}(x_1, x_n) \\ a_2 a_1 \text{Cov}(x_2, x_1) + a_2 a_2 \text{Cov}(x_2, x_2) + \dots + a_2 a_n \text{Cov}(x_2, x_n) \\ \vdots \\ a_n a_1 \text{Cov}(x_n, x_1) + a_n a_2 \text{Cov}(x_n, x_2) + \dots + a_n a_n \text{Cov}(x_n, x_n) \end{bmatrix}$$

reescrivendo por meio do produto de soma:

$$a_1 \begin{bmatrix} \text{Cov}(x_1, x_1) \\ \text{Cov}(x_2, x_1) \\ \vdots \\ \text{Cov}(x_n, x_1) \end{bmatrix} + a_2 \begin{bmatrix} \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_2) \\ \vdots \\ \text{Cov}(x_n, x_2) \end{bmatrix} + \dots + a_n \begin{bmatrix} \text{Cov}(x_1, x_n) \\ \text{Cov}(x_2, x_n) \\ \vdots \\ \text{Cov}(x_n, x_n) \end{bmatrix}$$

que continuando com a variável população
matricial temos:

$$\begin{bmatrix} \mu_1 & \mu_2 & \dots & \mu_n \end{bmatrix} \cdot \begin{bmatrix} \text{cov}(X_1, X_1) & \text{cov}(X_1, X_2) & \dots & \text{cov}(X_1, X_n) \\ \text{cov}(X_2, X_1) & \text{cov}(X_2, X_2) & \dots & \text{cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_n, X_1) & \text{cov}(X_n, X_2) & \dots & \text{cov}(X_n, X_n) \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

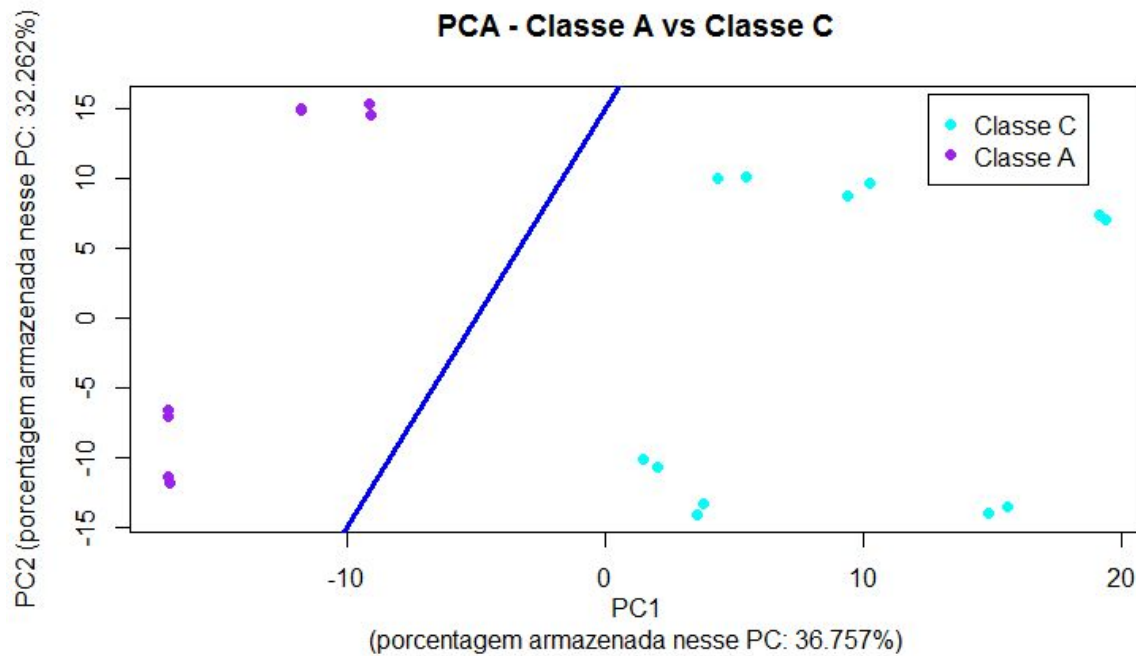
mat de cov

$\vec{\mu}_j^T$ $\vec{\mu}_j$

Então, $\vec{\mu}_j^T \sum \vec{\mu}_j$.

finalmente, $\text{Var}(\bar{z}_j) = \vec{\mu}_j^T \sum \vec{\mu}_j$

Questão 2



De acordo com o gráfico mostrado pelo código no arquivo zip e considerando a linha hipotética definida no gráfico em azul, as amostras tem agrupamento claro que pode ser visualizado e o plot possibilita essa visualização, se olharmos pra linha desenhada em azul e supondo que ela fosse a melhor fronteira de decisão para classificar uma nova amostra, fica claro que é possível usar os dois primeiros componentes principais para classificar as amostras nas classes descritas. As porcentagens das variâncias "armazenadas" no PC1 e no PC2 são 36.757% e 32.262% respectivamente.

Questão 3

O plot do PCA tridimensional pode ser uma enrascada tendo vista que o dependendo do referencial usado o gráfico não será informativo, por exemplo no caso da figura 4, é difícil observar as reais distâncias de uma amostra ou grupo de amostras para a outra amostra ou grupo de amostras. Além disso a imagem não mostra uma clara aproximação da as amostras de RNA amplificado de urina com as amostras tecido-específicas.

Uma questão é, que se o câncer analisado é de próstata, as amostras de urina vindas de indivíduos masculinos (provavelmente doentes) deveriam se aproximar de algum cubo azul (amostra característica) mostrando a similaridade, enquanto as amostras de urina de indivíduos femininos deveriam estar o mais longe possível dos cubos azuis mostrando que, de fato mulheres não poderiam ter câncer de próstata. Porém não fica claro essa diferenciação. Do ponto de vista de reconhecimento de padrões podemos ressaltar que não existe um plano que indique a separação dos clusters mencionados na descrição da figura.