

III 3 III

Amostragem multivariada

3.1. Introdução

Com os conceitos de álgebra vetorial introduzidos no capítulo 2, pode-se aprofundar na interpretação geométrica das estatísticas descritivas $\bar{\mathbf{X}}$, \mathbf{S} e \mathbf{R} . A maioria das explicações usam a representação das colunas de \mathbf{X} , como p pontos no espaço n dimensional. Será introduzida neste instante a pressuposição de que as observações constituem uma amostra aleatória. De uma forma simplificada, amostra aleatória significa (i) que as medidas tomadas em diferentes itens (unidades amostrais ou experimentais) são não relacionadas uma com as outras, e (ii) que a distribuição conjunta das p variáveis permanece a mesma para todos os itens. Essa estrutura de amostra aleatória é que justifica uma escolha particular de distância e dita a geometria para a representação n dimensional dos dados. Finalmente, quando os dados podem ser tratados como uma amostra aleatória à inferência estatística terá por base um sólido fundamento.

3.2. Geometria amostral

Uma observação multivariada é uma coleção de medidas em p variáveis tomadas na mesma unidade amostral ou experimental. No capítulo 1, item 1.3.1, as n observações obtidas foram dispostas em um arranjo (Matriz) \mathbf{X} por,

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2k} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{j1} & x_{j2} & \cdots & x_{jk} & \cdots & x_{jp} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} & \cdots & x_{np} \end{bmatrix}$$

em que cada linha de \mathbf{X} representa uma observação multivariada. Desde que o conjunto todo de mensurações é muitas vezes uma particular realização de variáveis aleatórias, diz-se que os dados representam uma amostra de tamanho n de uma população p variada.

Os dados podem ser plotados por um gráfico com p coordenadas. As colunas de \mathbf{X} representam n pontos no espaço p dimensional. Esse tipo de gráfico fornece informações de locação dos pontos e de variabilidade. Se os pontos pertencem a uma esfera, o vetor de médias amostrais, $\bar{\mathbf{X}}$, é o centro de balanço ou de massa. Se a variabilidade ocorre em mais de uma direção, pode-se detectar pela matriz de covariância, \mathbf{S} . Uma medida numérica única de variabilidade é fornecida pelo determinante da matriz de covariância.

Exemplo 3.1

Calcule o vetor média $\bar{\underline{x}}$ para a matriz \mathbf{X} apresentada a seguir. Plote os $n=3$ pontos no espaço $p=2$ (bidimensional) e localize $\bar{\underline{x}}$ no diagrama resultante.

$$\mathbf{X} = \begin{bmatrix} 2 & 1 \\ -3 & 0 \\ -2 & 2 \end{bmatrix}$$

A média amostral é dada por:

$$\bar{\underline{x}} = \begin{bmatrix} [(2+(-3))+(-2)]/3 \\ (1+0+2)/3 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

O primeiro ponto é dado por $\underline{x}'_1 = [2 \ 1]$, o segundo por $\underline{x}'_2 = [-3 \ 0]$, e o terceiro por $\underline{x}'_3 = [-2 \ 2]$. A Figura 3.1 mostra os pontos juntamente com $\bar{\underline{x}}$, centro de massa ou de balanço, obtidos pela matriz \mathbf{X} .

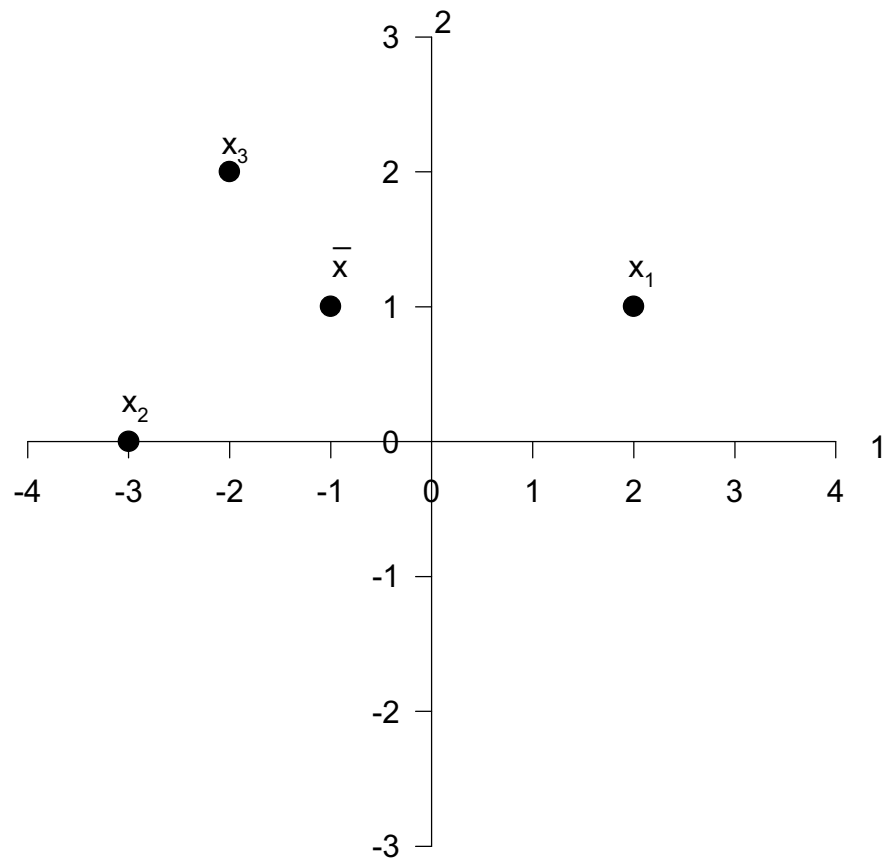


Figura 3.1. Diagrama com $n=3$ pontos no espaço bidimensional ($p=2$) mostrando o centro de massa, \bar{x} .

Uma representação alternativa é obtida através da consideração de p pontos no espaço n dimensional. Os elementos das linhas de \mathbf{X} são utilizados como coordenadas.

$$\begin{aligned}
 X &= \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2k} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{j1} & x_{j2} & \cdots & x_{jk} & \cdots & x_{jp} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} & \cdots & x_{np} \end{bmatrix} \\
 &= \begin{bmatrix} \tilde{y}_1 & \tilde{y}_2 & \cdots & \tilde{y}_k & \cdots & \tilde{y}_p \end{bmatrix}
 \end{aligned}$$

As coordenadas do k -ésimo ponto $\tilde{y}'_k = [x_{1k} \ x_{2k} \ \cdots \ x_{nk}]$ é determinada pela n -upla de todas as medidas da k -ésima variável. É conveniente representar \tilde{y}'_k como vetor ao invés de pontos.

Exemplo 3.2

Plote os dados da matriz X , com $p=2$ vetores no espaço tridimensional ($n=3$)

$$X = \begin{bmatrix} 2 & 1 \\ -3 & 0 \\ -3 & 2 \end{bmatrix}$$

$$\tilde{y}'_1 = [2 \ -3 \ -2] \quad \tilde{y}'_2 = [1 \ 0 \ 2]$$

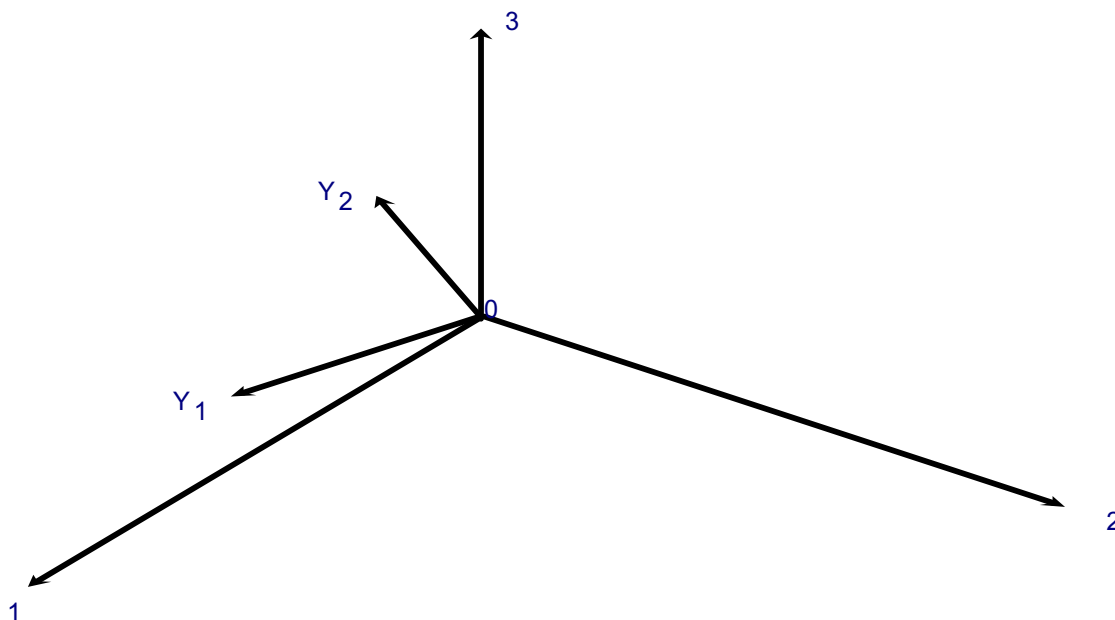


Figura 3.2. Diagrama da matriz de dados \mathbf{X} como $p=2$ vetores no espaço tridimensional.

Muita das expressões algébricas que serão encontradas na análise multivariada, podem ser relacionadas às noções geométricas de ângulos, comprimento (norma) e volumes. Isto é importante, pois representações geométricas facilitam a compreensão e conduzem a novas visões. Infelizmente, o ser humano está limitado a visualizar objetos no espaço tridimensional, e as representações da matriz \mathbf{X} não serão úteis se $n > 3$. No entanto, os relacionamentos geométricos e os conceitos estatísticos associados, descritos para o espaço tridimensional ou bidimensional, permanecem válidos para dimensões maiores.

É possível, em função do exposto, prover uma interpretação geométrica ao processo de encontrar a média amostral. O vetor $\underline{1}$ ($n \times 1$) será definido por $\underline{1}' = [1 \ 1 \ \dots \ 1]$. O vetor $\underline{1}$ forma um ângulo igual com cada um dos eixos coordenados,

de tal forma que $(1/\sqrt{n})\underline{1}$ tenha comprimento unitário e mesmo ângulo de direção.

Considerando o vetor $\underline{y}'_i = [x_{1i} \ x_{2i} \ \dots \ x_{ni}]$, cuja projeção em $(1/\sqrt{n})\underline{1}$ é:

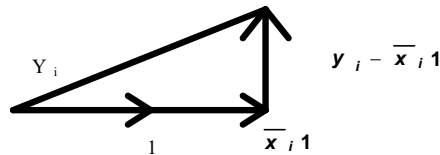
$$\underline{y}'_i \left(\frac{1}{\sqrt{n}} \underline{1} \right) \frac{1}{\sqrt{n}} \underline{1} = \frac{x_{1i} + x_{2i} + \dots + x_{ni}}{n} \underline{1} = \bar{x}_i \underline{1} = \underline{y}'_i \left(\frac{1}{n} \underline{1} \right) \underline{1}$$

Pois, a projeção geral de \underline{x} em \underline{y} é dada por:

$$\text{Proj}(\underline{X} \text{ em } \underline{Y}) = \frac{\underline{x}' \underline{y}}{|\underline{y}|} \underline{y}$$

Dessa forma $\bar{X}_i = \underline{y}'_i \left(\frac{1}{n} \underline{1} \right)$ corresponde a um múltiplo de 1, obtido a

partir da projeção de \underline{y}'_i em um vetor $\underline{1}$, de acordo com o esquema a seguir.



em que, $\underline{\tilde{y}}_i - \bar{x}_i \underline{1}$ é perpendicular a $\bar{x}_i \underline{1}$. Observe, também, que $\underline{e}_i = \underline{\tilde{y}}_i - \bar{x}_i \underline{1}$ é definido como desvios da i -ésima variável em relação a sua média amostral, e consiste nos elementos apresentados a seguir:

$$\underline{e}_i = \underline{\tilde{y}}_i - \bar{x}_i \underline{1} = \begin{bmatrix} \mathbf{X}_{i1} - \bar{x}_i \\ \mathbf{X}_{i2} - \bar{x}_i \\ \vdots \\ \mathbf{X}_{in} - \bar{x}_i \end{bmatrix},$$

A decomposição de \underline{y}_i , nos vetores média e desvio da média está apresentada esquematicamente na Figura 3.3 para $p=2$ e $n=3$.

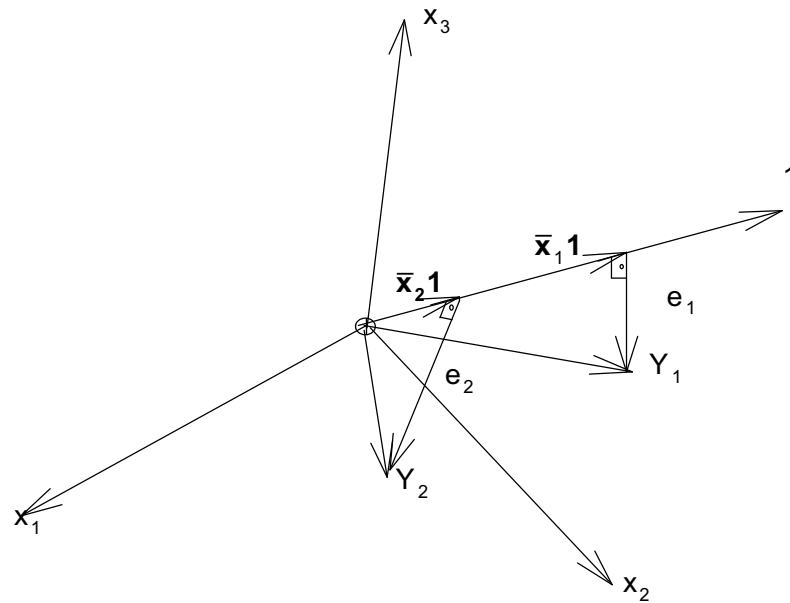


Figura 3.3. Decomposição de \underline{y}_i em componentes de média \bar{x}_{i1} e componentes de desvio $e_i = \underline{y}_i - \bar{x}_{i1}$.

Exemplo 3.3

Faça a decomposição de \underline{y}_i em componentes de média \bar{x}_{i1} e componentes de desvio $e_i = \underline{y}_i - \bar{x}_{i1}$, $i=1, 2$, para os dados do exemplo 3.2.

$$X = \begin{bmatrix} 2 & 1 \\ -3 & 0 \\ -3 & 2 \end{bmatrix} \quad \underset{\sim_1}{\mathbf{y}}' = [2 \quad -3 \quad -2] \quad \underset{\sim_2}{\mathbf{y}}' = [1 \quad 0 \quad 2]$$

$$\bar{x}_1 = \frac{2 + (-3) + (-2)}{3} = -1 \quad \bar{x}_2 = \frac{1 + 0 + 2}{3} = 1$$

$$\bar{x}_1 \underset{\sim}{\mathbf{1}} = -1 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} \quad \bar{x}_2 \underset{\sim}{\mathbf{1}} = 1 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\mathbf{e}_1 = \underset{\sim}{y}_1 - \bar{x}_1 \underset{\sim}{\mathbf{1}} = \begin{bmatrix} 2 \\ -3 \\ -2 \end{bmatrix} - \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} = \begin{bmatrix} 3 \\ -2 \\ -1 \end{bmatrix}$$

$$\mathbf{e}_2 = \underset{\sim}{y}_2 - \bar{x}_2 \underset{\sim}{\mathbf{1}} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}$$

Observa-se que: $\bar{x}_1 \underset{\sim}{\mathbf{1}}$ e \mathbf{e}_1 , $\bar{x}_1 \underset{\sim}{\mathbf{1}}$ e \mathbf{e}_2 , são perpendiculares.

$$(\bar{x}_1 \underset{\sim}{\mathbf{1}})' (\underset{\sim}{y}_1 - \bar{x}_1 \underset{\sim}{\mathbf{1}}) = [-1 \quad -1 \quad -1] \times \begin{bmatrix} 3 \\ -2 \\ -1 \end{bmatrix} = -3 + 2 + 1 = 0$$

A decomposição é:

$$\underset{\sim}{y}_1 = \begin{bmatrix} 2 \\ -3 \\ 2 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} + \begin{bmatrix} 3 \\ -2 \\ -1 \end{bmatrix}; \text{ e } \underset{\sim}{y}_2 = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}.$$

Os vetores de resíduos podem ser plotados a partir da origem, como apresentado na Figura 3.4, para os resíduos do exemplo 3.3.

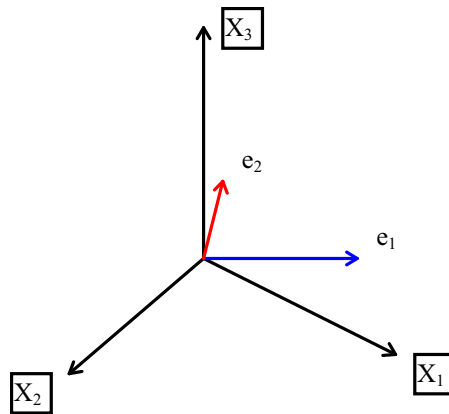


Figura 3.4. Vetores de desvios $\underset{\sim}{e}_i$ do exemplo 3.3.

Considere o comprimento ao quadrado dos vetores de desvios, obtidos por (2.2):

$$|\underset{\sim}{e}_i|^2 = \underset{\sim}{e}_i \cdot \underset{\sim}{e}_i = \sum_{j=1}^n (x_{ji} - \bar{x}_i)^2 \quad (3.1)$$

Observa-se por (3.1) que o comprimento ao quadrado dos vetores de desvios é proporcional à variância da i -ésima variável. Equivalentemente, o comprimento é proporcional ao desvio padrão. Vetores longos representam maiores variabilidades que os vetores mais curtos.

Para dois vetores desvios \underline{e}_i e \underline{e}_k :

$$\underline{e}_i' \underline{e}_k = \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) \quad (3.2)$$

De (2.3) e denotando o ângulo θ_{ik} como o ângulo formado pelos vetores \underline{e}_i e \underline{e}_k , tem-se:

$$\text{Cos}(\theta_{ik}) = \frac{\underline{e}_i' \underline{e}_k}{\sqrt{\underline{e}_i' \underline{e}_i} \times \sqrt{\underline{e}_k' \underline{e}_k}} \quad (3.3)$$

Usando (3.1) e (3.2) é fácil verificar que (3.3) é:

$$r_{ik} = \text{Cos}(\theta_{ik}) = \frac{S_{ik}}{\sqrt{S_{ii}} \sqrt{S_{kk}}} \quad (3.4)$$

O coseno do ângulo formado entre dois vetores desvios é igual ao coeficiente de correlação amostral. Portanto, se os dois vetores de desvios possuem a mesma orientação, o coeficiente de correlação será próximo de 1. Se os dois vetores

estão próximos de serem perpendiculares, a correlação amostral será próxima de zero. Se os dois vetores forem orientados em direções opostas, o coeficiente de correlação amostral será próximo de -1. Os conceitos de comprimento e ângulos permitem que se faça interpretações das estatísticas amostrais geometricamente, e auxiliam na compreensão dos seus significados.

3.3. Amostras aleatórias e esperanças do vetor média e da matriz de covariância amostral.

Com a finalidade de estudar a variabilidade amostral de estatísticas como \bar{X} e \mathbf{S} com a finalidade de se fazer inferências, é necessário fazer pressuposições a respeito das variáveis cujos valores observados constituem um conjunto de dados \mathbf{X} .

Supondo que os dados não foram ainda observados, mas pretende-se obter n mensurações em p variáveis. Antes de serem mensurados, os valores não podem em geral ser preditos exatamente. Conseqüentemente, estes são tratados como variáveis aleatórias. Neste contexto, os elementos (j, k) da matriz de dados representam realizações de uma variável aleatória, X_{jk} . Cada conjunto de medidas \underline{X}_j em p variáveis é um vetor aleatório.

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2k} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{j1} & x_{j2} & \cdots & x_{jk} & \cdots & x_{jp} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} & \cdots & x_{np} \end{bmatrix} = \begin{bmatrix} \underline{X}'_1 \\ \underline{X}'_2 \\ \vdots \\ \underline{X}'_j \\ \vdots \\ \underline{X}'_n \end{bmatrix} \quad (3.5)$$

Uma amostra aleatória pode ser definida por: “Se o vetor coluna $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$ em (3.5), representa independentes observações com distribuição conjunta com densidade $f(\underline{x})=f(x_1, x_2, \dots, x_p)$, então $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$ é uma amostra aleatória. Se a função conjunta de densidade é igual ao produto das marginais $f(\underline{x}_1) \cdot f(\underline{x}_2) \cdot \dots \cdot f(\underline{x}_n)$, sendo $f(\underline{x}_j)=f(x_{j1}, x_{j2}, \dots, x_{jp})$, então, $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$ é uma amostra aleatória.”

Algumas conclusões podem ser obtidas da distribuição de $\bar{\underline{X}}$ e \mathbf{S} sem pressuposições sobre a forma da distribuição conjunta das variáveis. Dessa forma, considere $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$ como uma amostra aleatória de uma distribuição conjunta com vetor média $\underline{\mu}$ e matriz de covariância Σ . então, $\bar{\underline{X}}$ é um estimador não viciado de $\underline{\mu}$ e sua matriz de covariância é $\frac{1}{n}\Sigma$. Isto é,

$$E(\bar{\underline{X}}) = \underline{\mu} \quad (\text{vetor média populacional})$$

$$\text{Cov}(\bar{\underline{X}}) = \frac{1}{n}\Sigma \quad (\text{Matriz de covariância populacional dividida pelo tamanho da amostra}).$$

PROVA:

$$\bar{\underline{X}} = (\underline{X}_1 + \underline{X}_2 + \dots + \underline{X}_n)/n$$

$$\begin{aligned}
E(\bar{X}) &= \left(\frac{1}{n}X_1 + \frac{1}{n}X_2 + \dots + \frac{1}{n}X_n\right) \\
&= E\left(\frac{1}{n}X_1\right) + E\left(\frac{1}{n}X_2\right) + \dots + E\left(\frac{1}{n}X_n\right) \\
&= \frac{1}{n} \left[nE(X_j) \right] = \frac{1}{n} \times n \times \mu
\end{aligned}$$

$$\therefore E(\bar{X}) = \mu$$

Para provar o valor da covariância, pode-se observar que:

$$(\bar{X} - \mu)(\bar{X} - \mu)' = \left(\frac{1}{n} \sum_{j=1}^n (X_j - \mu)\right) \times \left(\frac{1}{n} \sum_{\ell=1}^n (X_\ell - \mu)\right)' = \frac{1}{n^2} \sum_{j=1}^n \sum_{\ell=1}^n (X_j - \mu)(X_\ell - \mu)'$$

Então,

$$\text{Cov}(\bar{X}) = E(\bar{X} - \mu)(\bar{X} - \mu)' = \frac{1}{n^2} \sum_{j=1}^n \sum_{\ell=1}^n E(X_j - \mu)(X_\ell - \mu)'$$

Sendo $j \neq \ell$ e considerando que $E(X_j - \mu)(X_\ell - \mu)'$ é igual a zero, devido a covariância entre os elementos independentes X_j e X_ℓ ser nula, então,

$$\text{Cov}(\bar{X}) = \frac{1}{n^2} \sum_{j=1}^n E(X_j - \mu)(X_j - \mu)'$$

Desde que $\Sigma = E(\tilde{X}_j - \underline{\mu})(\tilde{X}_j - \underline{\mu})'$ é a covariância populacional comum dos componentes \tilde{X}_j , têm-se:

$$\begin{aligned} \text{Cov}(\bar{\tilde{X}}) &= \frac{1}{n^2} \sum_{j=1}^n E(\tilde{X}_j - \underline{\mu})(\tilde{X}_j - \underline{\mu})' = \frac{1}{n^2} (\Sigma + \Sigma + \dots + \Sigma) = \\ &= \frac{1}{n^2} (n\Sigma) = \frac{1}{n} \Sigma \end{aligned}$$

3.4. Variância Generalizada

Com uma única variável, a variância da amostra é usada para descrever a variação nas mensurações desta variável. Quando p variáveis são observadas em cada unidade da amostra ou do experimento, a variação é descrita pela matriz de variância e covariância amostral.

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} & \dots & \mathbf{S}_{1p} \\ \mathbf{S}_{21} & \mathbf{S}_{22} & \dots & \mathbf{S}_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{S}_{p1} & \mathbf{S}_{p2} & \dots & \mathbf{S}_{pp} \end{bmatrix}$$

A matriz de covariância amostral contém p variâncias e $\frac{1}{2}p(p-1)$ covariâncias, potencialmente diferentes. Algumas vezes, no entanto, deseja-se expressar a variação por um único valor numérico. Uma escolha deste valor é o

determinante de \mathbf{S} , o qual reduz à variância amostral usual para o caso de uma única variável ($p=1$). Este determinante é denominado de variância amostral generalizada.

$$\text{Variância amostral Generalizada} = |\mathbf{S}| \quad (3.6)$$

Exemplo 3.4

O peso de espiga PE (\mathbf{X}_1), e o número de espigas NE (\mathbf{X}_2), foi avaliado em 28 variedades de milho em Sete Lagoas, MG. A matriz de covariância amostral \mathbf{S} , obtida dos dados é:

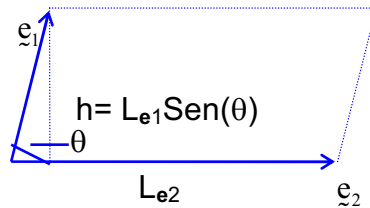
$$\mathbf{S} = \begin{bmatrix} 2,905 & 9,096 \\ 9,096 & 90,817 \end{bmatrix}$$

A variância generalizada neste caso é:

$$\text{Variância amostral Generalizada} = |\mathbf{S}| = 2,905 \times 90,817 - 9,096^2 = 181,0862$$

A variância amostral generalizada se constitui numa forma de escrever toda a informação de todas as variâncias e covariâncias como um único valor numérico. Obviamente, quando $p > 1$ é possível que algumas informações amostrais sejam perdidas no processo. A interpretação geométrica, no entanto, poderá mostrar a força e as fraquezas desta estatística descritiva.

Considerando-se o volume (área) gerado no plano definido por dois vetores de desvios $\underline{\epsilon}_1 = Y_1 - \bar{X}_1 \underline{1}$ e $\underline{\epsilon}_2 = Y_2 - \bar{X}_2 \underline{1}$. Seja L_{ϵ_1} e L_{ϵ_2} os comprimentos dos vetores $\underline{\epsilon}_1$ e $\underline{\epsilon}_2$, respectivamente. Da geometria têm-se:



A área do triângulo é $L_{\epsilon_1} \times \text{Sen}(\theta) \times L_{\epsilon_2}$, podendo ser expressa por:

$$\text{Área} = L_{\epsilon_1} L_{\epsilon_2} \sqrt{1 - \cos^2(\theta)}$$

Mas,

$$L_{\epsilon_1} = \sqrt{\sum_{j=1}^n (X_{j1} - \bar{X}_1)^2} = \sqrt{(n-1)S_{11}}$$

$$L_{\epsilon_2} = \sqrt{\sum_{j=1}^n (X_{j2} - \bar{X}_2)^2} = \sqrt{(n-1)S_{22}}$$

$$\text{Cos}(\theta) = r_{12}$$

Portanto,

$$\text{Área} = (n-1)\sqrt{S_{11}S_{22}(1-r_{12}^2)} \quad (3.7)$$

Por outro lado,

$$\begin{aligned} |\mathbf{S}| &= \begin{vmatrix} S_{11} & S_{21} \\ S_{12} & S_{22} \end{vmatrix} = \begin{vmatrix} S_{11} & \sqrt{S_{11}}\sqrt{S_{22}}r_{12} \\ \sqrt{S_{11}}\sqrt{S_{22}}r_{12} & S_{22} \end{vmatrix} \\ &= S_{11}S_{22} - S_{11}S_{22}r_{12}^2 = S_{11}S_{22}(1-r_{12}^2) \end{aligned} \quad (3.8)$$

Se (3.7) e (3.8) forem comparados, pode-se observar que

$$|\mathbf{S}| = (\text{Área})^2 / (n-1)^2$$

Esta expressão pode ser generalizada para p vetores desvios por indução:

$$\text{Variância amostral Generalizada} = |\mathbf{S}| = (\text{Volume})^2 \cdot (n-1)^{-p} \quad (3.9)$$

A equação (3.9) mostra que a variância amostral é proporcional ao quadrado do volume gerado pelos p vetores desvios. A Figura 3.5 (a) e (b) mostra as regiões trapezoidais geradas com $p=3$ vetores resíduos, correspondentes a “grandes” e “pequenas” variâncias amostrais generalizadas, respectivamente.

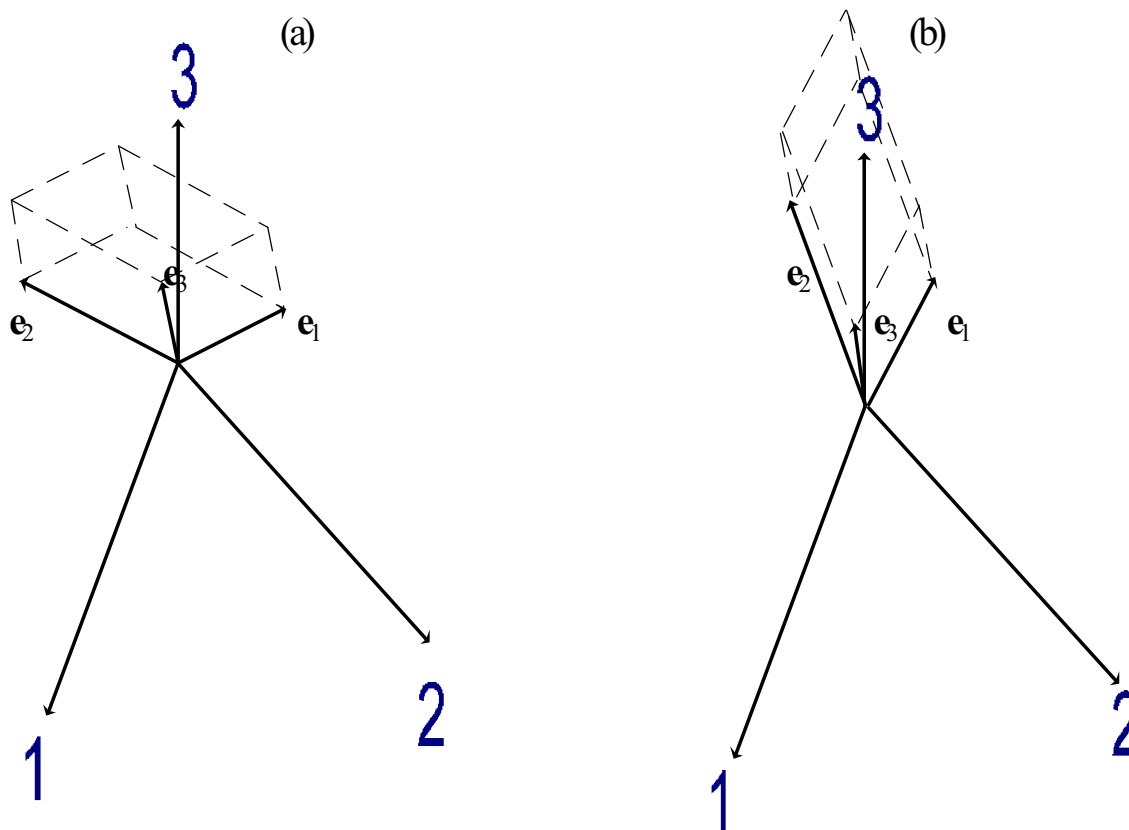


Figura 3.5. (a) grande variância amostral generalizada, e (b) pequena variância amostral generalizada, para $p=3$.

Para um tamanho amostral fixo, é óbvio que $|S|$ cresce com o aumento do comprimento dos vetores de desvios e_i (ou $\sqrt{(n-1)S_{ii}}$). Em adição, o volume aumentará para um comprimento fixado, se os vetores residuais forem movidos até possuírem ângulos retos. Por outro lado se um ou mais dos vetores residuais aproximarem do hiperplano formado por outros vetores residuais, o volume diminuirá tendendo a zero.

Apesar da variância amostral generalizada possuir algumas interpretações geométricas formidáveis como as ilustradas na Figura 3.5, ela sofre alguns problemas

como estatística amostral capaz de sumarizar a informação contida na matriz \mathbf{S} . Para ilustrar estas deficiências, considere as matrizes de covariâncias e os coeficientes de correlações apresentados a seguir.

$$\mathbf{S} = \begin{bmatrix} 10 & 8 \\ 8 & 10 \end{bmatrix} \quad \mathbf{S} = \begin{bmatrix} 10 & -8 \\ -8 & 10 \end{bmatrix} \quad \mathbf{S} = \begin{bmatrix} 6 & 0 \\ 0 & 6 \end{bmatrix}$$

$$r_{12} = \frac{8}{\sqrt{10}\sqrt{10}} = 0,8 \quad r_{12} = \frac{-8}{\sqrt{10}\sqrt{10}} = -0,8 \quad r_{12} = \frac{0}{\sqrt{6}\sqrt{6}} = 0,0$$

$$|\mathbf{S}| = 36$$

$$|\mathbf{S}| = 36$$

$$|\mathbf{S}| = 36$$

Apesar das três matrizes possuírem a mesma variância amostral generalizada ($|\mathbf{S}|=36$), elas possuem estruturas de correlações distintas. Portanto, diferentes estruturas de correlações não são detectadas pela variância amostral generalizada. As situações em que $p>2$ podem ser ainda mais obscuras.

Muitas vezes é desejável mais informações do que um simples valor como $|\mathbf{S}|$ pode oferecer como resumo de \mathbf{S} . Pode-se mostrar que $|\mathbf{S}|$ pode ser expresso como produto dos autovalores de \mathbf{S} ($|\mathbf{S}|=\lambda_1.\lambda_2....\lambda_p$). A elipsóide centrada na média é baseada em \mathbf{S}^{-1} , possui eixos de comprimento proporcionais a raiz quadrada de λ_i 's de \mathbf{S} , que reflete a variabilidade no sentido do i -ésimo autovalor. Esta elipsóide é apresentada a seguir.

$$(\underline{\mathbf{X}} - \bar{\underline{\mathbf{X}}})' \mathbf{S}^{-1} (\underline{\mathbf{X}} - \bar{\underline{\mathbf{X}}}) = c^2 \quad (3.10)$$

Demonstra-se que o volume desta hiperelipsóide é proporcional à raiz quadrada de $|\mathbf{S}|$. Desta forma, os autovalores, fornecem informações da variabilidade em todas as direções da representação no espaço p-dimensional dos dados. Portanto, é mais útil apresentar seus valores individuais do que seu produto. Este tópico será abordado com mais detalhe quando se discutir sobre os componentes principais.

A variância amostral generalizada será zero se um dos vetores residuais pertencer a um (hiper) plano formado por uma combinação linear dos outros, ou seja, quando as linhas da matriz de desvios, forem linearmente dependentes.

Exemplo 3.5

Mostre que $|\mathbf{S}|=0$ para

$$X = \begin{bmatrix} 3 & 3 & 6 \\ 1 & 3 & 4 \\ 2 & 0 & 2 \end{bmatrix}$$

O vetor média é:

$$\bar{X}' = [2 \quad 2 \quad 4]$$

Os vetores dos desvios são:

$$X - 1\bar{X}' = [e_1 \quad e_2 \quad e_3] = \begin{bmatrix} 1 & 1 & 2 \\ -1 & 1 & 0 \\ 0 & -2 & -2 \end{bmatrix}$$

Verifica-se que $\underline{e}'_3 = \underline{e}'_1 + \underline{e}'_2$, ou seja:

$$[2 \ 0 \ -2] = [1 \ -1 \ 0] + [1 \ 1 \ -2] = [2 \ 0 \ -2] \text{ c.q.d.}$$

Isto significa que um dos vetores resíduos, pertence ao plano gerado pelos outros dois. Desta forma o volume tridimensional é zero (degenerescência). Este caso é ilustrado na Figura 3.6 e demonstrado numericamente através da obtenção de $|\mathbf{S}|$.

$$\mathbf{S} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 3 & 3 \\ 1 & 3 & 4 \end{bmatrix}$$

Pela definição (2.9), têm-se:

$$\begin{aligned} |\mathbf{S}| &= 1 \times \begin{vmatrix} 3 & 3 \\ 3 & 4 \end{vmatrix} \times (-1)^2 + 0 \times \begin{vmatrix} 0 & 1 \\ 3 & 4 \end{vmatrix} \times (-1)^3 + 1 \times \begin{vmatrix} 0 & 1 \\ 3 & 3 \end{vmatrix} \times (-1)^4 = \\ &= 1.3.1 + 0 + 1.(-3).1 = 3 - 3 = 0 \end{aligned}$$

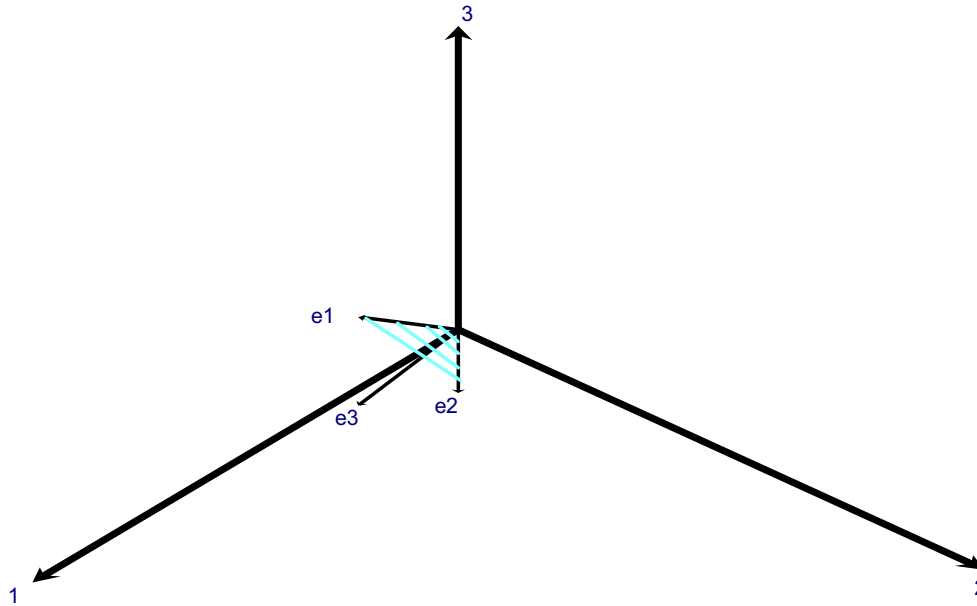


Figura 3.6 Caso em que $|\mathbf{S}|=0$ (degenerescência) para o volume tridimensional.

Em qualquer análise estatística o resultado $|\mathbf{S}|=0$ indica que existem variáveis redundantes, ou seja, que possuem a mesma informação, e que estas podem ser removidas do estudo. A matriz de covariância reduzida, será de posto completo e a variância generalizada diferente de zero. A questão de quais variáveis remover no caso de degenerescência não é fácil de responder e será abordado nos estudos de componentes principais. No entanto, quando há possibilidade de escolha, o pesquisador deve reter as medidas de uma variável (presumidamente) causal ao invés de uma característica secundária.

3.5.Variância generalizada de variáveis padronizadas

A variância amostral generalizada é influenciada pela diferença de variabilidade das mensurações das variáveis individuais, ou seja, caso a variância amostral de uma determinada variável (S_{ii}) seja grande ou pequena em relação às demais. O vetor residual correspondente $\underline{e}_i = \underline{Y}_i - \bar{x}_i \underline{1}$ será muito longo ou muito curto, do ponto de vista geométrico e terá um papel importante na determinação do volume. É muitas vezes necessário, em função do exposto, padronizar os vetores residuais, de tal forma que eles tenham o mesmo comprimento.

A padronização destes vetores residuais é equivalente a transformar as variáveis originais x_{jk} pelos seus valores $(x_{jk} - \bar{x}_k) / \sqrt{S_{kk}}$. A matriz de covariância amostral das variáveis padronizadas será então igual a \mathbf{R} , ou seja, igual a matriz de correlação das variáveis originais. Dessa forma pode-se definir:

$$\text{Variância generalizada amostral das variáveis padronizadas} = |\mathbf{R}| \quad (3.11)$$

Os vetores resíduos resultantes, cujos valores são dados por $\mathbf{e}_{jk} = (x_{jk} - \bar{x}_k) / \sqrt{S_{kk}}$, possuem todos comprimento igual a $\sqrt{n-1}$. A variância generalizada amostral das variáveis padronizadas será grande se estes vetores forem perpendiculares e pequena se dois ou mais deles tiverem próximas da mesma direção. Em (3.4) foi visto que o coseno do ângulo θ_{ik} entre os vetores residuais \underline{e}_i e \underline{e}_k , com $i \neq k$, é igual ao coeficiente de correlação amostral r_{ik} . Dessa forma, o $|\mathbf{R}|$ será grande quando todos os r_{ik} forem próximos de zero e será pequeno quando um ou mais dos r_{ik} for próximo de -1 ou de +1.

Utilizando os mesmos argumentos que conduziram a (3.9) pode-se verificar que:

$$|\mathbf{R}| = (n-1)^{-p} (\text{volume})^2 \quad (3.12)$$

O volume gerado pelos vetores desvios de $p=3$ variáveis padronizadas está ilustrado na Figura 3.7. Estes vetores desvios padronizados são correspondentes aos vetores desvios da Figura 3.5, cuja comparação revela que a influência do vetor \underline{e}_2 (com grande variabilidade na direção de x_2) no volume quadrado de $|\mathbf{S}|$ é maior do que sua influência no volume quadrado de $|\mathbf{R}|$.

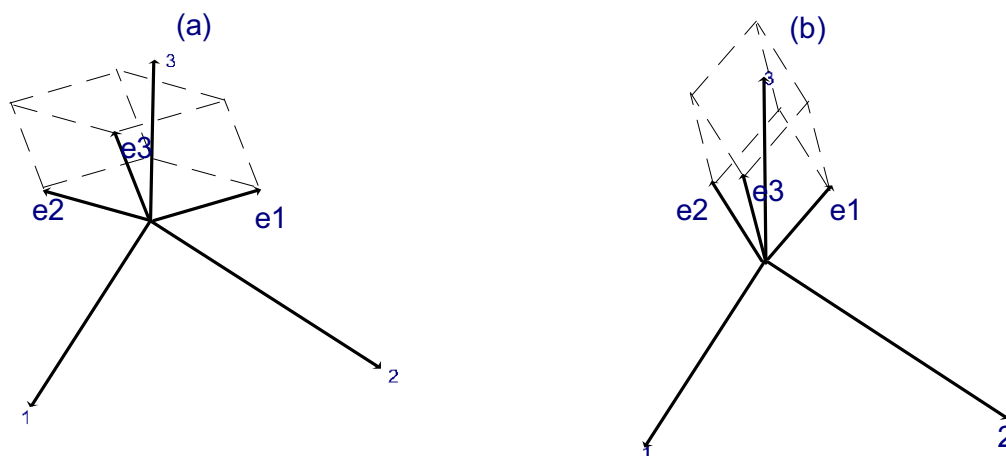


Figura 3.7. Volume gerado por três variáveis padronizadas: (a) grande variância e (b) pequena variância generalizada.

As quantidades $|\mathbf{S}|$ e $|\mathbf{R}|$ são relacionadas por:

$$|\mathbf{S}| = (S_{11} S_{22} \dots S_{pp}) |\mathbf{R}| \quad (3.13)$$

Exemplo 3.6

É ilustrada através deste exemplo a relação (3.13) entre $|\mathbf{S}|$ e $|\mathbf{R}|$ para $p=3$ caracteres de milho (x_1 : diâmetro do colmo; x_2 : número de folhas; e x_3 : comprimento de folhas). A matriz \mathbf{R} e \mathbf{S} obtidas são:

$$\mathbf{S} = \begin{bmatrix} 4,935 & 0,552 & 2,921 \\ 0,552 & 0,686 & 1,932 \\ 2,921 & 1,932 & 17,993 \end{bmatrix} \text{ e } \mathbf{R} = \begin{bmatrix} 1,00 & 0,30 & 0,31 \\ 0,30 & 1,00 & 0,55 \\ 0,31 & 0,55 & 1,00 \end{bmatrix}$$

Usando-se a definição de determinante (2.9), tem-se:

$$|\mathbf{S}|=37,3878$$

$$|\mathbf{R}|=0,6137$$

Usando (3.13) e os resultados obtidos:

$$|\mathbf{S}| = (S_{11} S_{22} \dots S_{pp}) |\mathbf{R}|$$

$$37,3878 = (4,935 \times 0,686 \times 17,993) \times 0,6137$$

$37,3878 \approx 37,3828$ (verificado, apesar da pequena diferença devido às aproximações nos cálculos)

3.6. Outra generalização da variância

Uma outra medida capaz de sintetizar a informação contida na matriz de covariância que é utilizada em componentes principais é definida pela soma dos elementos da diagonal da matriz de covariância \mathbf{S} e é denominada de variância amostral total. Portanto,

$$\text{Variância amostral total} = \text{Traço de } \mathbf{S} = \text{Tr}(\mathbf{S}) = S_{11} + S_{22} + \dots + S_{pp}$$

(3.14)

Exemplo 3.7

Calcular a variância amostral total da matriz \mathbf{S} do exemplo (3.6)

$$\text{Tr}(\mathbf{S}) = S_{11} + S_{22} + S_{33} = 4,935 + 0,686 + 17,993 = 23,614$$

Geometricamente a variância amostral total representa a soma dos comprimentos ao quadrado dos vetores residuais \mathbf{e}_i ($i=1, 2, \dots, p$) dividido por $n-1$. Ela não considera as orientações dos vetores residuais, sendo portanto limitada para ser

utilizada com variáveis padronizadas, pois seu valor será sempre o mesmo para distintos conjuntos de dados desde que o número de variáveis destes seja igual.

3.7. Exercícios

3.7.1. Plote os $n=4$ pontos no diagrama bidimensional e localize \bar{X} no diagrama resultante.

$$X = \begin{bmatrix} 1 & 1 \\ -1 & -1 \\ -1 & 1 \\ 1 & -1 \end{bmatrix}$$

3.7.2. Encontre o ângulo entre os vetores \underline{y}_1 e \underline{y}_2 do exemplo 3.1. Calcule o coseno do mesmo e discuta sobre o significado deste resultado.

3.7.3. Obtenha a decomposição dos vetores \underline{y}_1 e \underline{y}_2 do exemplo 3.1 em componente de média e componente de desvio. Comprove a ortogonalidade dos componentes de média com os vetores de desvios ou residuais.

3.7.4. Calcule usando (3.3) o coseno do ângulo entre os vetores residuais \underline{e}_1 e \underline{e}_2 obtidos em 3.3. Calcule o coeficiente de correlação usando (1.4) entre as variáveis 1 e 2, e compare os resultados obtidos.

3.7.5. Obtenha as matrizes de covariância amostral para o conjunto de dados do exercício 3.7.1, e calcule as variâncias amostrais generalizadas das variáveis originais e padronizadas. Calcule também a variância amostral total.

3.7.6. Qual é a área do trapezóide gerado pelos $p=2$ vetores desvios, do exercício 3.7.1.

3.8. Referências

BOCK, R.D. **Multivariate statistical methods in behavioral research**. McGraw-Hill, 1975.

JOHNSON, R.A.; WICHERN, D.W. **Applied multivariate statistical analysis**. 4th edition. Prentice Hall, New Jersey, 1998. 816p.