

1.3. THE DATA MATRIX

$$X_{n \times p} = \begin{bmatrix} x_{n1} \\ \vdots \\ x_{np} \end{bmatrix} = [x_{n(1)} \dots x_{n(p)}] = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix}$$

1.4. Summary statistics

$$\bar{x}_i = \frac{1}{n} \sum_{n=1}^N x_{ni} \quad \text{SAMPLE MEAN (variable i)}$$

$$\bar{x} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix} \quad \text{SAMPLE MEAN VECTOR}$$

$$\text{for } \bar{x}_{n \times 1} = \begin{bmatrix} x_{n1} \\ \vdots \\ x_{np} \end{bmatrix}$$

$$= \frac{1}{n} \sum_{n=1}^N \bar{x}_{n \times 1} = \frac{1}{n} X^T \mathbf{1}_n$$

$$= \frac{1}{n} \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix}_{p \times n} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1}$$

$$= \frac{1}{n} \left(\sum_{n=1}^N x_{n1} + \dots + \sum_{n=1}^N x_{np} \right)$$

$$\sigma_{ii} = \frac{1}{n} \sum_{n=1}^n (x_{ni} - \bar{x}_i)^2 = \sigma_i^2 \quad \text{SAMPLE VARIANCE}$$

$$\sigma_{ij} = \frac{1}{n} \sum_{n=1}^n (x_{ni} - \bar{x}_i)(x_{nj} - \bar{x}_j) \quad \text{SAMPLE COVARIANCE}$$

$$= \frac{1}{n} \sum_{n=1}^n (x_{ni}x_{nj} - \cancel{x_{ni}}\bar{x}_j - \bar{x}_i\cancel{x_{nj}} + \bar{x}_i\bar{x}_j)$$

$$\Sigma = \begin{bmatrix} \sigma_{11} & \dots & \sigma_{1p} \\ \vdots & \ddots & \vdots \\ \sigma_{n1} & \dots & \sigma_{np} \end{bmatrix} \quad \text{SAMPLE VARIANCE-COVARIANCE MATRIX}$$

↳ possui efeitos de escala

$$\Sigma = \frac{1}{n} \sum_{n=1}^n (x_{nn} - \bar{x}_n)(x_{nn} - \bar{x}_n)^T$$

$$= \frac{1}{n} X^T X - \bar{x}_n \bar{x}_n^T = \frac{1}{n} (X^T X - \frac{1}{n} X^T \frac{1}{n} \frac{1}{n}^T X)$$

$$X^T X = \begin{bmatrix} \sum_{n=1}^n x_{n1}^2 & \dots & \sum_{n=1}^n x_{n1} x_{np} \\ \vdots & & \vdots \\ \sum_{n=1}^n x_{np} x_{np} & \dots & \sum_{n=1}^n x_{np}^2 \end{bmatrix}_{p \times p}$$

$$\frac{1}{n} \frac{1^T}{n} = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}_{n \times n}$$

$$X^T \frac{1}{n} \frac{1^T}{n} = \begin{bmatrix} \bar{x}_{n1} & \dots & \bar{x}_{nL} \\ \vdots & \ddots & \vdots \\ \bar{x}_{np} & \dots & \bar{x}_{np} \end{bmatrix}_{p \times n}$$

$$\begin{aligned}
X^T \frac{1}{n} H^T X &= \begin{bmatrix} \bar{x}_{n_1} & \dots & \bar{x}_{n_p} \\ \vdots & \ddots & \vdots \\ \bar{x}_{n_p} & \dots & \bar{x}_{n_1} \end{bmatrix} \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix} \\
&= \begin{bmatrix} (\bar{x}_{n_1})^2 & \dots & (\bar{x}_{n_1})(\bar{x}_{n_p}) \\ \vdots & \ddots & \vdots \\ (\bar{x}_{n_p})(\bar{x}_{n_1}) & \dots & (\bar{x}_{n_p})^2 \end{bmatrix} \\
&= \frac{1}{n} X^T H X \quad (\text{Quadratic form})
\end{aligned}$$

$$H = I - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \quad \text{CENTRING matrix}$$

$$= \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & & \vdots \\ \vdots & \ddots & \ddots & 1 \\ 0 & \dots & \dots & 1 \end{bmatrix} - \begin{bmatrix} \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \\ \frac{1}{n} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \end{bmatrix} = \begin{bmatrix} \frac{n-1}{n} & -\frac{1}{n} & \dots & -\frac{1}{n} \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ -\frac{1}{n} & \dots & \dots & \frac{n-1}{n} \end{bmatrix}$$

$H = H^T \Rightarrow H$ is **idempotent**

- Since H is idempotent, it follows that, for any p -vector α

$$\underline{\alpha}^T S \underline{\alpha} = \frac{1}{n} \underline{\alpha}^T X^T H^T H X \underline{\alpha} = \frac{1}{n} \underline{y}^T \underline{y} \geq 0$$

$$\underline{y} = H X \underline{\alpha} \quad \text{to prove that if } n > p \Rightarrow S \geq 0$$

$\therefore S \geq 0$ (**positive semi-definite**)

$$M = \sum_{n=1}^N \underline{x}_n \underline{x}_n^T = X^T X$$

o distância de manhattan
 matrix of sums of squares and products

$$r_{ij} = s_{ij} / (\sqrt{s_i s_j})$$

$$R = (r_{ij}) \geq 0$$

$$= D^{-1} S D^{-1}$$

$$D = \text{DIAG}(s_i) \quad \text{Am}$$

distância de Pearson
correlation coefficient

distância euclidiana (retângular)

sample correlation matrix
 É "paranormalizada" a covariância

1. Pearson
2. Spearman
3. correl. Pearsais
4. matriz de importâncias
5. copulas
6. distância entre observações

$$\text{trace}(S) = \text{trace}(\Lambda) = \text{trace}(\text{diag}(\lambda_i)) = \sum_{i=1}^p \lambda_i$$

Variancia Total

$$(\text{b}) |\Delta| = \prod_{i=1}^p \lambda_i$$

Variancia Generalizada

1.5 LINEAR COMBINATIONS

- Consider the following linear combination:

$$y_n = a_1 x_{n1} + \dots + a_p x_{np} \quad n=1, \dots, N$$

where $a_i, i=1, \dots, p$, are given.

- The mean of y is given by:

$$\bar{y} = \frac{1}{N} \sum_{n=1}^N a_n^T \bar{x}_n = a^T \bar{x}$$

$$\begin{bmatrix} a_1 & \dots & a_p \end{bmatrix} \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix} = \sum_{i=1}^p a_i \bar{x}_i$$

linear combination
of the mean of
each column

and the variance

$$\sigma^2_y = \frac{1}{N} \sum_{n=1}^N (y_n - \bar{y})^2 = \frac{1}{N} \sum_{n=1}^N a_n^T (x_n - \bar{x})(x_n - \bar{x})^T$$

- In general, we may be interested in q -dimensional ($q \leq p$) linear transformation

$$y_n = A x_n + b, \quad n=1, \dots, N$$

$$Y = X A^T + b^T, \quad A_{q \times p}, \quad b_{q \times 1}$$

- THE MEAN VECTOR AND COVARIANCE MATRIX FOR \tilde{y}_n IS GIVEN BY :

$$\bar{\tilde{y}}_n = A\bar{x}_n + \bar{b}$$

$$S_{\tilde{y}} = \frac{1}{n} \sum_{n=1}^n (\tilde{y}_n - \bar{\tilde{y}}_n)(\tilde{y}_n - \bar{\tilde{y}}_n)^T = ASA^T$$

- IF A IS NON-INVERTIBLE, THEN

$$S = A^{-1} S_{\tilde{y}} (A^T)^{-1}$$

- SOME EXAMPLES OF LINEAR TRANSFORMATIONS ARE :

1. Scaling Transformation: SCALES EACH VARIABLE TO HAVE UNIT VARIANCE

$$* y_n = D^{-1} (x_n - \bar{x}_n) , \quad n = 1, \dots, n$$

$$D = \text{diag}(s_i)$$

2. Mahalanobis Transformation: ELIMINATES THE CORRELATION BETWEEN VARIABLES AND STANDARDIZE ITS VARIANCE.

if $S > 0 \Rightarrow S^{-1}$ HAS UNIQUE POSITIVE SQUARE ROOT $S^{-1/2}$ AND

$$* z_n = S^{-1/2} (x_n - \bar{x}_n) , \quad n = 1, \dots, n$$

3. Principal Component transformation: By the Spectral decomposition theorem the covariance matrix may be written as:

COVARIANCE MATRIX

spectral decomposition

$$S = G L G^T, \quad G G^T = G^T G = I \text{ (idempotent)}$$

G : An orthogonal matrix of eigenvectors

L : Diagonal matrix of the eigen values of S
 $(\lambda_1 > \lambda_2 > \dots > \lambda_p > 0)$

$$\therefore \underline{w}_n = G^T (\underline{x}_n - \bar{\underline{x}})$$

SINCE $S_w = G^T S G = L$ is diagonal, the columns of W , called principal components, represent uncorrelated linear combination of the variables.

→ characteristic roots of S (eigenvalues)

$$|S - \lambda I_p| = 0 \Rightarrow \lambda = [\lambda_1, \lambda_2, \dots, \lambda_p]$$

$$\Rightarrow L = \text{DIAG}(\lambda)$$

Exercise

$$S = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix}, \quad S_{11} = S_{22}$$

Compute the eigenvalues and eigenvectors

- Note**
- If $m > p \Rightarrow S$ HAS $\text{RANK}(S) \leq p \Rightarrow$
 \Rightarrow at most p non-zero EIGENVALUES
 - If $m \leq p \Rightarrow S$ HAS $m-p$ null EIGENVALUES

Exercise

state the relationships between:

1. M AND p dimensions of a matrix S
2. $\text{RANK}(S)$
3. eigenvalues of S

u

ESTATÍSTICAS DAS UNIDADES AMOSTRAIS

ESTATÍSTICAS DESCRIPTIVAS DAS
UNIDADES AMOSTRAIS

R-técnicas
Q-técnicas

$D_{n \times n} = (d_{i,j}^2)$ DISTÂNCIA EUCLIDIANA AMOSTRAL

$$d_{i,j}^2 = d_{i,j}^2 = (x_i - x_j)^T (x_i - x_j)$$

: $\mathbb{R}^p \rightarrow \mathbb{R}$

$$d_{Pik}^2 = (x_i^* - x_k^*)^T (x_i^* - x_k^*) = (x_i - \bar{x}_k)^T D_{S_{jj}}^{-1} (x_i - \bar{x}_k)$$

↳ DISTÂNCIA EUCLIDIANA PADRONIZADA (PEARSON)

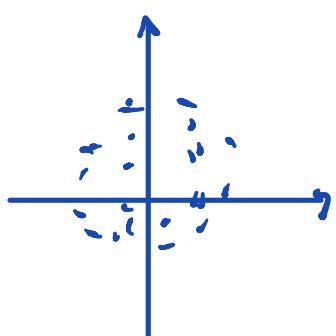
$$= 2(1 - \rho_{ik})$$

↳ COEFICIENTE DE CORRELAÇÃO DE PEARSON

$$d_{Min}^2 = z_i^T z_i = (y_i - \bar{y})^T S^{-1} (y_i - \bar{y})$$

DISTÂNCIA DE
MHALANOBIS

Interpretação Geométrica



VARIÁVEIS:

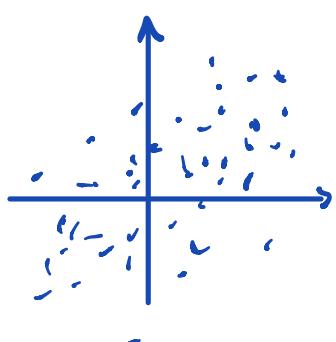
a) INDEPENDENTES $d_E^2 = (x - \bar{x})^T (x - \bar{x})$

b) HOMOCOVARÍATAS

a) INDEPENDENTES

b) HETEROCOVARÍATAS

$$d_P^2 = (x - \bar{x})^T D_{S_{jj}}^{-1} (x - \bar{x})$$



a) DEPENDENTES

b) HETEROCOVARÍATAS

$$d_H^2 = (y - \bar{y})^T S^{-1} (y - \bar{y})$$