

## Técnicas Multivariadas em Saúde

Lupércio França Bessegato  
Dep. Estatística/UFJF

Técnicas Multivariadas em Saúde - 2015

## Roteiro

1. Introdução
2. Distribuições de Probabilidade Multivariadas
3. Representação de Dados Multivariados
4. Testes de Significância  $c/$  Dados Multivariados
5. Análise de Componentes Principais
6. Análise Fatorial
7. Análise de Correlação Canônica
8. Análise de Conglomerados
9. Análise Discriminante
10. Análise de Correspondência
11. Referências

Técnicas Multivariadas em Saúde - 2015

## Análise de Componentes Principais

Técnicas Multivariadas em Saúde - 2015

## Introdução

- Objetivo:
  - √ Explicar a estrutura de variância e covariância de conjunto de variáveis através de algumas combinações lineares das mesmas
  - √ Busca-se:
    - Redução de dados
    - Interpretação

Técnicas Multivariadas em Saúde - 2015

### Componentes Principais Exatas

- Algebricamente:
  - √ Combinações lineares particulares das  $p$  variáveis aleatórias  $X_1, X_2, \dots, X_p$ .
- Geometricamente:
  - √ Representam a seleção de um novo sistema de coordenadas obtidas por rotação do sistema original
  - √ Os novos eixos representam as direções com maior variabilidade
  - √ Fornecem descrição mais simples e mais parcimoniosa da estrutura de covariâncias

Técnicas Multivariadas em Saúde - 2015

- Componentes principais:

- √ São necessárias  $p$  componentes para reproduzir a variabilidade total do sistema
- √ As componentes são não correlacionadas entre si
  - Ortogonalidade entre as componentes
- √ Variabilidade das  $p$  variáveis é aproximada pela variabilidade das  $k$  principais componentes
  - Buscam-se situações em que haja quase tanta informação nas  $k$  componentes principais quanto nas  $p$  variáveis originais

Técnicas Multivariadas em Saúde - 2015

- Análise de componentes principais:

- √ Não pressupõe normalidade
  - Componentes principais derivadas de populações normais têm interpretações úteis
- √ Com frequência, revela relações insuspeitadas
  - Pode permitir interpretações que não seriam obtidas preliminarmente
- √ Em geral, é um passo intermediário para a aplicação de outras técnicas

Técnicas Multivariadas em Saúde - 2015

### Componentes Principais Exatas Extraídas da Matriz de Covariâncias

- Sejam o vetor aleatório

$$\mathbf{X}' = [X_1, X_2, \dots, X_p].$$

com matriz de covariâncias é  $\Sigma$ , cujos autovalores são  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ .

- Componentes principais de  $\Sigma$ :

$$Y_1, Y_2, \dots, Y_p.$$

- √ Combinações lineares não correlacionadas do vetor aleatório, cujas variâncias são as maiores possíveis

Técnicas Multivariadas em Saúde - 2015

- Definição 1 – Componente principal:

√ A j-ésima componente principal da matriz  $\Sigma$  é definida como:

$$Y_j = \mathbf{e}'_j \mathbf{X} = e_{j1}X_1 + e_{j2}X_2 + \dots + e_{jp}X_p.$$

√  $\mathbf{e}_j$ : autovetor correspondente ao j-ésimo autovalor

- Esperança e variância de  $Y_j$ :

$$E[Y_j] = E[\mathbf{e}'_j \mathbf{X}] = \mathbf{e}'_j \boldsymbol{\mu} = e_{j1}\mu_1 + e_{j2}\mu_2 + \dots + e_{jp}\mu_p.$$

$$\text{Var}[Y_j] = \text{Var}[\mathbf{e}'_j \mathbf{X}] = \mathbf{e}'_j \Sigma \mathbf{e}_j = \mathbf{e}'_j \left( \sum_{i=1}^p \mathbf{e}_i \mathbf{e}'_i \right) \mathbf{e}_j = \lambda_j.$$

- Covariância entre duas componentes principais:

$$\text{Cov}[Y_j, Y_k] = 0, j \neq k$$

Técnicas Multivariadas em Saúde - 2015

- Comentário:

√ Cada autovalor  $\lambda_j$  representa a variância de uma componente principal  $Y_j$ .

√ Autovalores estão ordenados em ordem decrescente

– A primeira componente é a de maior variabilidade

– A p-ésima componente é a de menor variabilidade

Técnicas Multivariadas em Saúde - 2015

- Variâncias total e generalizada de  $\Sigma$ :

√ Total:  $\text{tr}(\Sigma) = \sum_{i=1}^p \sigma_{ii} = \sum_{i=1}^p \lambda_i$

√ Generalizada de  $\Sigma$ :  $|\Sigma| = \prod_{i=1}^p \lambda_i$

√ Em termos dessas duas medidas globais de variação, os vetores  $\mathbf{X}$  e  $\mathbf{Y}$  são equivalentes

Técnicas Multivariadas em Saúde - 2015

- Proporção da variância total que é explicada pela j-ésima componente principal:

$$\frac{\text{Var}[Y_j]}{\text{Variância total de } \mathbf{X}} = \frac{\lambda_j}{\text{tr}(\Sigma)} = \frac{\lambda_j}{\sum_{i=1}^p \lambda_i}$$

√ 1ª componente tem a maior proporção de explicação

- Proporção da variância total que é explicada pelas k primeiras componentes principais

$$\frac{\sum_{j=1}^k \text{Var}[Y_j]}{\text{Variância total de } \mathbf{X}} = \frac{\sum_{j=1}^k \lambda_j}{\text{tr}(\Sigma)} = \frac{\sum_{j=1}^k \lambda_j}{\sum_{i=1}^p \lambda_i}$$

√ Busca-se analisar um conjunto menor de variáveis sem perder muita informação sobre a estrutura de variabilidade original

Técnicas Multivariadas em Saúde - 2015

• Aproximação de  $\Sigma$ :

- √ Analisando as k primeiras componentes principais

$$\Sigma_{p \times p} \approx \sum_{i=1}^k \lambda_i \mathbf{e}_i \mathbf{e}_i'$$

- √ Cada parcela da soma envolve uma matriz de dimensão p x p correspondente apenas à informação da j-ésima componente principal

Técnicas Multivariadas em Saúde - 2015

• Definição 2 – Componente principal:

- √ Sistema cuja j-ésima combinação linear de  $\mathbf{X}$  é definida como:

$$Y_j = \mathbf{a}_j' \mathbf{X} = a_{j1}X_1 + a_{j2}X_2 + \dots + a_{jp}X_p.$$

- √  $\mathbf{e}_j$ : autovetor correspondente ao j-ésimo autovalor

• Esperança e variância de  $Y_j$ :

$$E[Y_j] = E[\mathbf{e}_j' \mathbf{X}] = \mathbf{e}_j' \boldsymbol{\mu} = e_{j1}\mu_1 + e_{j2}\mu_2 + \dots + e_{jp}\mu_p.$$

$$\text{Var}[Y_j] = \text{Var}[\mathbf{a}_j' \mathbf{X}] = \mathbf{a}_j' \boldsymbol{\Sigma} \mathbf{a}_j.$$

• Covariância entre duas componentes principais:

$$\text{Cov}[Y_j, Y_k] = \mathbf{a}_j' \boldsymbol{\Sigma} \mathbf{a}_k, \quad j \neq k, \quad j = 1, 2, \dots, p$$

Técnicas Multivariadas em Saúde - 2015

- √ Buscam-se os valores dos coeficientes  $a_{ij}$ , tais que:

- i.  $Y_1, Y_2, \dots, Y_p$  tenham variância máxima e sejam não correlacionadas entre si
- ii. Os vetores  $\mathbf{a}_i$  tenham comprimento unitário:

$$\mathbf{a}_j' \mathbf{a}_k = \begin{cases} 1 & , \text{ se } j = k \\ 0 & , \text{ se } j \neq k \end{cases}$$

- √ Pode-se demonstrar que :

- A variância máxima de  $(\mathbf{a}_i' \mathbf{X})$  é igual a  $\lambda_i$ .
- É obtida quando  $\mathbf{a}_i = \mathbf{e}_i$ .

Técnicas Multivariadas em Saúde - 2015

**Correlação entre Componente Principal e Variável Aleatória**

- Os coeficientes de correlação entre a componente principal  $Y_i$  de S e a variável  $X_k$  é

$$\rho_{Y_i, X_k} = \frac{e_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}$$

- √ A magnitude de  $e_{ik}$  mede a contribuição da k-ésima variável na i-ésima componente (a despeito das outras variáveis).

- Não medem a importância de  $X_k$  na presença das outras variáveis.
- Alguns estatísticos recomendam que somente os valores  $e_{ik}$  (e não as correlações) sejam consideradas na interpretação dos componentes

Técnicas Multivariadas em Saúde - 2015

### Estimação das Componentes Principais – Matriz de Covariâncias

- Em geral,  $\Sigma$  é estimada por  $S$ :

$$S = \begin{bmatrix} S_{11} & S_{12} & \dots & S_{1p} \\ S_{12} & S_{22} & \dots & S_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ S_{1p} & S_{2p} & \dots & S_{pp} \end{bmatrix}$$

√ Autovalores de  $S$ :  $\hat{\lambda}_1, \hat{\lambda}_2, \dots, \hat{\lambda}_p$

√ Autovetores de  $S$ :  $\hat{e}_1, \hat{e}_2, \dots, \hat{e}_p$

Técnicas Multivariadas em Saúde - 2015

- $j$ -ésima componente principal de  $S$ :

$$\hat{Y}_j = \hat{e}'_j \mathbf{X} = \hat{e}_{j1}X_1 + \hat{e}_{j2}X_2 + \dots + \hat{e}_{jp}X_p, \quad j = 1, 2, \dots, p.$$

- Componentes principais amostrais – Propriedades

- Variância:  $\text{Var}[\hat{Y}_j] = \hat{\lambda}_j$ .
- Covariância entre as componentes:  $\text{Cov}(\hat{Y}_j, \hat{Y}_k) = 0, \quad j \neq k$
- Variância total estimada explicada pela componente:

$$\frac{\text{Var}[\hat{Y}_j]}{\text{Variância total estimada de } \mathbf{X}} = \frac{\hat{\lambda}_j}{\text{tr}(\mathbf{S})} = \frac{\hat{\lambda}_j}{\sum_{i=1}^p \hat{\lambda}_i}$$

- Correlação estimada entre componente e variável:

$$r_{\hat{Y}_j, X_k} = \frac{\hat{e}_{jk} \sqrt{\hat{\lambda}_j}}{\sqrt{S_{kk}}}$$

Técnicas Multivariadas em Saúde - 2015

- Decomposição espectral de  $S$ :

$$S = \sum_{j=1}^p \hat{\lambda}_j \mathbf{e}_j \mathbf{e}'_j$$

√ Aproximação de  $S$  pelas primeiras  $k$  componentes

$$S_{p \times p} \approx \sum_{i=1}^k \hat{\lambda}_i \hat{\mathbf{e}}_i \hat{\mathbf{e}}'_i$$

- Scores das componentes

√ Valor das componentes para cada elemento amostral

√ Na prática, o uso das componentes relevantes dá através dos scores

Técnicas Multivariadas em Saúde - 2015

### Exemplo 8.3

- Pesquisa com 5 variáveis socioeconômicas

√  $X_1$ : população total (milhares)

√  $X_2$ : Escolaridade mediana (anos concluídos)

√  $X_3$ : Emprego total (milhares)

√  $X_4$ : Empregos na área da saúde (centenas)

√  $X_5$ : Valor mediano da habitação (x \$10.000)

- Dados: *BD\_multivariada.xls/pesquisa*

Técnicas Multivariadas em Saúde - 2015

- Vetor de médias amostral ( $\bar{x}$ )

Variable	Mean
X1_Pop	4,323
X2_escol	14,014
X3_empregos	1,952
X4_saude	2,171
X5_habitacao	2,454

- Matriz de covariâncias amostral (S)

Covariances: X1\_Pop; X2\_escol; X3\_empregos; X4\_saude; X5\_habitacao

	X1_Pop	X2_escol	X3_empregos	X4_saude	X5_habitacao
X1_Pop	4,307556				
X2_escol	1,683680	1,767473			
X3_empregos	1,802776	0,588026	0,800669		
X4_saude	2,155326	0,177978	1,064828	1,969475	
X5_habitacao	-0,253474	0,175549	-0,158339	-0,356807	0,504380

- A variação amostral pode ser resumida por uma ou duas componentes principais?

Técnicas Multivariadas em Saúde - 2015

- Correlação mede unicamente importância de uma variável individual sem considerar a influência das demais

✓ No exemplo, os coeficientes de correlação confirmam a interpretação fornecida pelos coeficientes das componentes

Técnicas Multivariadas em Saúde - 2015

	Componentes Principais						
	1	2	3	4	5		
	e1	r(y1,xk)	e2	r(y2,xk)	e3	e4	e5
População Total	0,781	0,99	0,071	-0,04	-0,004	-0,542	0,302
Escolaridade Mediana	0,306	0,61	0,764	-0,76	0,162	0,545	0,009
Total de Empregos	0,334	0,98	-0,083	0,12	-0,015	-0,051	-0,937
Empregos Área Saúde	0,426	0,80	-0,579	0,55	-0,220	0,636	0,172
Valor Mediano Habitação	-0,054	-0,20	0,262	0,49	-0,962	-0,051	-0,025
Variância	6,931		1,785		0,390	0,230	0,014
% Variância Total (acumulada)	74,1		93,2		97,4	99,8	100,0

- Variância amostral é bem resumida por 2 componentes
  - ✓ redução de 14 observações de 5 variáveis para 14 observações de 2 variáveis
  - ✓ 1ª. componente: média ponderada de 4 variáveis
  - ✓ 2ª. componente: contraste entre empregos saúde com média ponderada da escolaridade com valor habitação

Técnicas Multivariadas em Saúde - 2015

### Número de Componentes Principais

- Quantas componentes principais devem ser retidas?
  - ✓ Não há resposta definitiva
- Considerações a serem tomadas:
  - ✓ Quantidade explicada de variância amostral total
  - ✓ Tamanho relativo dos autovalores (variância das componentes amostrais)
  - ✓ Interpretação das componentes

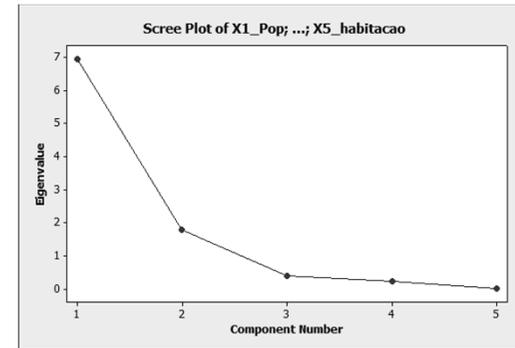
Técnicas Multivariadas em Saúde - 2015

### Scree Plot

- Gráfico  $\lambda_i$  vs.  $i$ 
  - √ Procura-se um 'cotovelo' no gráfico
  - √ São consideradas as componentes até o ponto em que os autovalores remanescentes são relativamente pequenos e todos aproximadamente do mesmo valor

Técnicas Multivariadas em Saúde - 2015

### Exemplo 8.3



Técnicas Multivariadas em Saúde - 2015

### Exemplo 8.4

- Relação entre tamanho e forma de cascos de tartaruga
  - √ Comprimento
  - √ Largura
  - √ Espessura
  - √ Gênero: male/female
- Análise para as tartarugas macho
- Literatura sugere transformação logarítmica em estudos de relação entre tamanho e forma
- Dados: *BD\_multivariada.xls/tartarugas*

Técnicas Multivariadas em Saúde - 2015

### Vetor de médias amostral ( $\bar{x}$ )

Variable	Mean
log_comp_male	4,7254
log_larg_male	4,4776
log_esp_male	3,7032

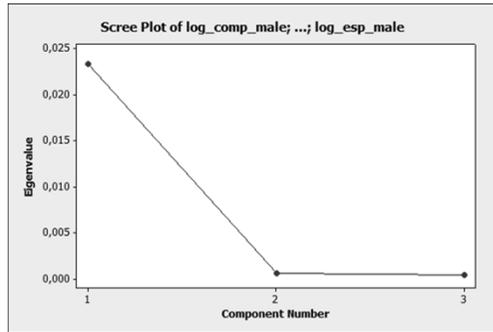
### Matriz de covariâncias amostral (S)

	log_comp_male	log_larg_male	log_esp_male
log_comp_male	0,01107200		
log_larg_male	0,00801914	0,00641673	
log_esp_male	0,00815965	0,00600527	0,00677276

- A variação amostral pode ser resumida por uma principal?

Técnicas Multivariadas em Saúde - 2015

• Scree Plot



√ Uma componente principal é claramente dominante

Técnicas Multivariadas em Saúde - 2015

• Componentes principais:

Principal Component Analysis: log\_comp\_male; log\_larg\_male; log\_esp\_male

Eigenanalysis of the Covariance Matrix

Eigenvalue	0,023303	0,000598	0,000360
Proportion	0,961	0,025	0,015
Cumulative	0,961	0,985	1,000

Variable	PC1	PC2	PC3
log_comp_male	0,683	-0,159	-0,713
log_larg_male	0,510	-0,594	0,622
log_esp_male	0,523	0,788	0,324

• Componente adotada:

$$\hat{y}_1 = 0,683 \ln(comp) + 0,510 \ln(larg) + 0,523 \ln(espes)$$

$$= \ln [(comp)^{0,683} (larg)^{0,510} (esp)^{0,523}]$$

√ ln(volume) de uma caixa com dimensões ajustadas

Técnicas Multivariadas em Saúde - 2015

**Componentes Principais de Variáveis Padronizadas**

• Padronização do vetor aleatório  $\mathbf{X}$ :

$$\mathbf{Z} = (\mathbf{V}^{1/2})^{-1} (\mathbf{X} - \boldsymbol{\mu})$$

√  $\mathbf{V}^{1/2}$ : matriz diagonal de desvios-padrão

√ Variável padronizada:  $Z_i = \frac{X_i - \mu_i}{\sqrt{\sigma_{ii}}}$

√ Matriz de covariâncias de  $\mathbf{Z}$ :

$$\text{Cov}(\mathbf{Z}) = (\mathbf{V}^{1/2})^{-1} \boldsymbol{\Sigma} (\mathbf{V}^{1/2})^{-1} = \mathbf{P}$$

√ Componentes principais de  $\mathbf{Z}$ :

- Obtidas dos autovalores e autovetores de  $\mathbf{P}$ .

Técnicas Multivariadas em Saúde - 2015

• Componente principal das variáveis padronizadas:

√ A j-ésima componente principal da matriz  $\boldsymbol{\Sigma}$ :

$$Y_j = \mathbf{e}'_j \mathbf{Z} = \mathbf{e}'_j (\mathbf{V}^{1/2})^{-1} (\mathbf{X} - \boldsymbol{\mu}) = e_{j1} Z_1 + e_{j2} Z_2 + \dots + e_{jp} Z_p$$

√  $\mathbf{e}_j$ : autovetor da matriz de correlações  $\mathbf{P}$ .

• Variância total de  $\mathbf{P}$ :

$$\sum_{j=1}^p \text{Var}[Y_j] = \sum_{j=1}^p \text{Var}[Z_j] = p$$

√ Proporção de variância populacional (padronizada) devido à j-ésima componente

$$\frac{\text{Var}[Y_j]}{\text{Variância total de } \mathbf{Z}} = \frac{\lambda_j}{\text{tr}(\mathbf{P})} = \frac{\lambda_j}{p}, k = 1, 2, \dots, p$$

√ Correlação entre  $Y_j$  e  $X_k$ :  $\rho_{Y_j, X_k} = e_{jk} \sqrt{\lambda_j}, i, k = 1, 2, \dots, p$

Técnicas Multivariadas em Saúde - 2015

### Comentários

- As componentes principais de  $\Sigma$  são diferentes daquelas obtidas de  $\mathbf{P}$ .
  - √ Seus autovalores e autovetores são diferentes
  - √ Um conjunto de componentes principais não é simplesmente uma função do outro conjunto
- A padronização traz consequências
  - √ Variáveis deveriam ser padronizadas se elas são medidas em escalas com amplitudes muito diferentes
    - Ex. Vendas anuais e razão entre lucro/ativos

Técnicas Multivariadas em Saúde - 2015

### Padronização dos Componentes Principais Amostrais

- Frequentemente são padronizadas:
  - √ Variáveis medidas em diferentes escalas
  - √ Na mesma escala, mas com amplitudes bastante diferentes
- As componentes principais não são invariantes às mudanças na escala

Técnicas Multivariadas em Saúde - 2015

- Padronização dos elementos amostrais:

$$\mathbf{z}_j = \mathbf{D}^{-1/2} (\mathbf{x}_j - \bar{\mathbf{x}}) = \begin{bmatrix} \frac{x_{j1} - \bar{x}_1}{\sqrt{s_{11}}} \\ \frac{x_{j2} - \bar{x}_2}{\sqrt{s_{22}}} \\ \vdots \\ \frac{x_{jp} - \bar{x}_p}{\sqrt{s_{pp}}} \end{bmatrix}, j = 1, 2, \dots, n$$

√  $\mathbf{D}$ : matriz diagonal dos desvios-padrão amostrais

- Matriz de dados:

$$\mathbf{Z}_{n \times p} = \begin{bmatrix} z_{11} & z_{12} & \dots & z_{1p} \\ z_{21} & z_{22} & \dots & z_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ z_{n1} & z_{n2} & \dots & z_{np} \end{bmatrix} = \begin{bmatrix} \mathbf{z}'_1 \\ \mathbf{z}'_2 \\ \vdots \\ \mathbf{z}'_n \end{bmatrix}.$$

Técnicas Multivariadas em Saúde - 2015

### Análise de Componentes Principais – Matriz de Correlações

- As componentes principais obtidas a partir da matriz de covariâncias são influenciadas pelas variáveis de maior variância
  - √ A padronização das variáveis ameniza esse problema
- Análise de componentes principais de variáveis padronizadas é equivalente a obter as componentes principais através da matriz de correlações

Técnicas Multivariadas em Saúde - 2015

### Estimação das Componentes Principais – Matriz de Correlação

- **P** é estimada por **R**:

√ Importante:  $S_Z = R$

$$R = \begin{bmatrix} 1 & r_{12} & \dots & r_{1p} \\ r_{12} & 1 & \dots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ r_{1p} & r_{2p} & \dots & 1 \end{bmatrix}$$

√ Autovalores de **R**:  $\hat{\lambda}_1, \hat{\lambda}_2, \dots, \hat{\lambda}_p$

√ Autovetores de **S**:  $\hat{e}_1, \hat{e}_2, \dots, \hat{e}_p$

Técnicas Multivariadas em Saúde - 2015

- **j**-ésima componente principal de **R**:

$$\hat{Y}_j = \hat{e}'_j Z = \hat{e}_{j1} Z_1 + \hat{e}_{j2} Z_2 + \dots + \hat{e}_{jp} Z_p, \quad j = 1, 2, \dots, p.$$

- Componentes principais amostrais – Propriedades

i. Variância:  $\text{Var}[\hat{Y}_j] = \hat{\lambda}_j$ .

ii. Covariância entre as componentes:  $\text{Cov}(\hat{Y}_j, \hat{Y}_k) = 0, \quad j \neq k$

iii. Variância total estimada explicada pela componente:

$$\frac{\hat{\lambda}_j}{p}$$

iv. Correlação estimada entre componente e variável:

$$r_{\hat{Y}_j, X_k} = \frac{\hat{e}_{jk} \sqrt{\hat{\lambda}_j}}{\sqrt{S_{kk}}}$$

Técnicas Multivariadas em Saúde - 2015

### Exemplo 8.5

- Taxas de retorno de 5 ações negociadas na Bolsa de New York

√ Período: Jan./75 a Dez./76

√ Ações:

- Allied Chemical
- du Pont
- Union Carbide
- Exxon
- Texaco

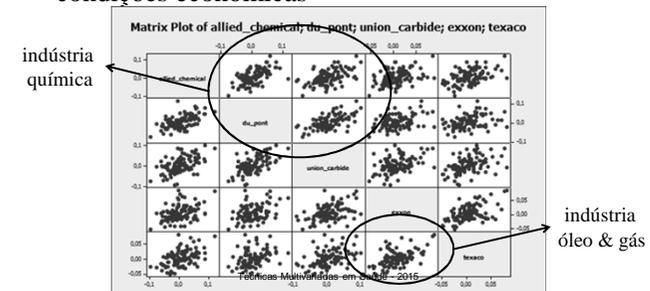
√ Dados: BD\_multivariada.xls/

Técnicas Multivariadas em Saúde - 2015

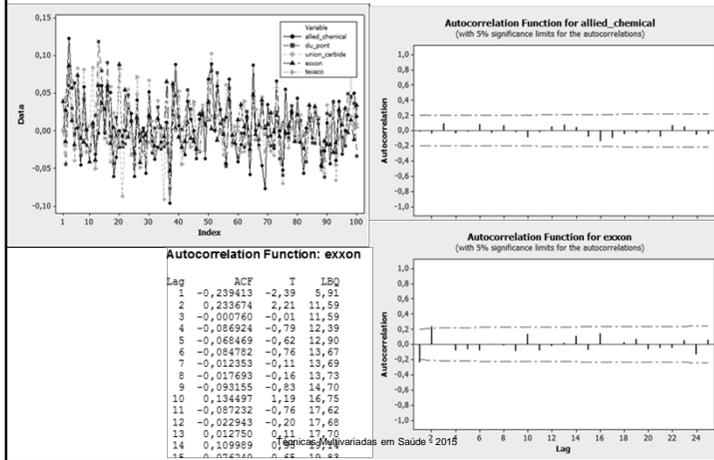
- Taxa de retorno:  $(\text{taxa de retorno}) = \frac{(\$ \text{ fechamento sexta atual}) - (\$ \text{ fechamento sexta anterior})}{(\$ \text{ fechamento sexta anterior})}$

- As taxas de retorno entre ativos estão correlacionadas

√ ações tendem a se mover juntas em resposta às condições econômicas



- Observações de 100 semanas aparentam estar distribuídas independentemente



LFB1

- Vetor de médias amostral ( $\bar{x}$ )

**Descriptive Statistics: allied\_chemical; du\_pont; union\_carbide; Exxon; texaco**

Variable	Mean
allied_chemical	0,00543
du_pont	0,00483
union_carbide	0,00565
Exxon	0,00629
texaco	0,00371

- Matriz de correlação amostral (S)

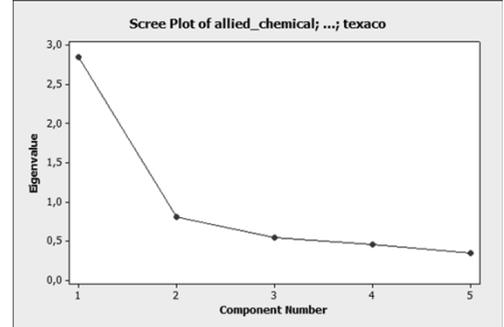
**Correlations: allied\_chemical; du\_pont; union\_carbide; Exxon; texaco**

	allied_chemical	du_pont	union_carbide	Exxon
du_pont	0,577			
union_carbide	0,509	0,598		
Exxon	0,387	0,390	0,436	
texaco	0,462	0,322	0,426	0,524

- A variação amostral pode ser resumida por uma ou duas componentes principais?

Técnicas Multivariadas em Saúde - 2015

- Scree Plot



√ Aparentemente duas componentes principais resumem bem os dados

Técnicas Multivariadas em Saúde - 2015

- Componentes principais:

**Principal Component Analysis: allied\_chemi; du\_pont; union\_carbid; Exxon; texac**

**Eigenanalysis of the Correlation Matrix**

Eigenvalue	2,8565	0,8091	0,5400	0,4513	0,3430
Proportion	0,571	0,162	0,108	0,090	0,069
Cumulative	0,571	0,733	0,841	0,931	1,000

Variable	PC1	PC2	PC3	PC4	PC5
allied_chemical	0,464	0,241	0,613	-0,381	-0,453
du_pont	0,457	0,509	-0,178	-0,211	0,675
union_carbide	0,470	0,261	-0,337	0,664	-0,396
Exxon	0,422	-0,525	-0,539	-0,473	-0,179
texaco	0,421	-0,582	0,434	0,381	0,387

√ Duas primeiras componentes com 73% da variabilidade amostral padronizada total

Técnicas Multivariadas em Saúde - 2015

**Slide 56**

---

**LFB1**

Calcular matriz de covariâncias amostral

Há domínio de variabilidade?

Lupércio Bessegato; 20/02/2013

• 1ª. componente principal:

$$\hat{y}_1 = 0,464z_1 + 0,457z_2 + 0,470z_3 + 0,421z_4 + 0,421z_5$$

√ Variáveis:

- $z_1$ : retorno padronizado – Allied Chemical
- $z_1$ : retorno padronizado – du Pont
- $z_1$ : retorno padronizado – Union Carbide
- $z_1$ : retorno padronizado – Exxon
- $z_1$ : retorno padronizado – Texaco

√ Interpretação:

- soma ponderada (índice) das 5 ações
- pesos aproximadamente iguais
- Componente geral do mercado de ações  
(componente do mercado)

Técnicas Multivariadas em Saúde - 2015

• 2ª. componente principal:

$$\hat{y}_1 = 0,240z_1 + 0,509z_2 + 0,260z_3 - 0,526z_4 - 0,582z_5$$

√ Interpretação:

- contraste entre ações de indústrias químicas e de óleo & gás
- Componente industrial

Técnicas Multivariadas em Saúde - 2015

• Comentários:

√ A maioria das variações dos ativos devem-se às atividades de mercado (1ª. componente) e atividades industriais não correlacionadas (2ª. componente)

√ As componente remanescentes não são de simples interpretação

- coletivamente, representam variação que é provavelmente específica de cada ação

Técnicas Multivariadas em Saúde - 2015

### Variáveis Padronizadas – Regra Empírica

- Reter apenas as componentes cujas variâncias ( $\lambda_i$ ) são maiores que a unidade
  - √ componente que explicam individualmente pelo menos  $1/p$  da variância amostral padronizada total
- No caso do exemplo anterior (8.6), pareceu-se sensível reter uma componente ( $y_2$ ) associada à autovalor menor que a unidade

Técnicas Multivariadas em Saúde - 2015

### Importante

- √ Um valor pequeno incomum para o último autovalor da matriz de covariâncias (ou correlação) amostral pode indicar uma dependência linear não detectada no conjunto de dados
- √ Valores grande de autovalores (e correspondentes autovetores são importantes em uma análise
- √ Autovalores próximos de zero não devem ser ignorados
  - Autovetores associados podem apontar dependências lineares no conjunto de dados (problemas computacionais ou de interpretação)

Técnicas Multivariadas em Saúde - 2015

### Gráfico dos Componentes Principais

- Podem:
  - √ revelar observações suspeitas
  - √ fornecer verificações da hipótese de normalidade

Técnicas Multivariadas em Saúde - 2015

- São combinações das variáveis originais:
  - √ Se as observações provém de população normal multivariada, é razoável esperar que as componentes sejam aproximadamente normais
  - √ Se forem usadas como entrada em análises adicionais
    - Verificar se as 1<sup>a</sup>.s componentes são aproximadamente normais
- As últimas componentes principais podem ajudar a apontar observações suspeitas

Técnicas Multivariadas em Saúde - 2015

### Resumo

- Procedimento auxiliar na verificação de normalidade
  - √ Construir diagrama de dispersão para os pares dos primeiros componentes principais
  - √ Construir Q-Q plots para os valores amostrais gerados por cada componente principal
- Identificação de observações suspeitas:
  - √ Construir diagramas de dispersão e Q-Q plots para as últimas componentes principais.

Técnicas Multivariadas em Saúde - 2015

### Exemplo 8.7

- Plotando os Componentes Principais dos dados das tartarugas macho:

$$\sqrt{x_1} = \ln(\text{comp})$$

$$\sqrt{x_2} = \ln(\text{larg})$$

$$\sqrt{x_3} = \ln(\text{esp})$$

- Componentes:

$$\hat{y}_1 = 0,683 \ln(x_1 - 4,725) + 0,510 \ln(x_2 - 4,478) + 0,523 \ln(x_3 - 3,703)$$

$$\hat{y}_2 = -0,159 \ln(x_1 - 4,725) - 0,594 \ln(x_2 - 4,478) + 0,788 \ln(x_3 - 3,703)$$

$$\hat{y}_3 = -0,713 \ln(x_1 - 4,725) + 0,622 \ln(x_2 - 4,478) + 0,324 \ln(x_3 - 3,703)$$

Técnicas Multivariadas em Saúde - 2015

- Comandos Minitab para Q-Q Plot

```
Name C30 "(j-1/2)/n"
Set C30
1(1 : 24 / 1) 1
End.

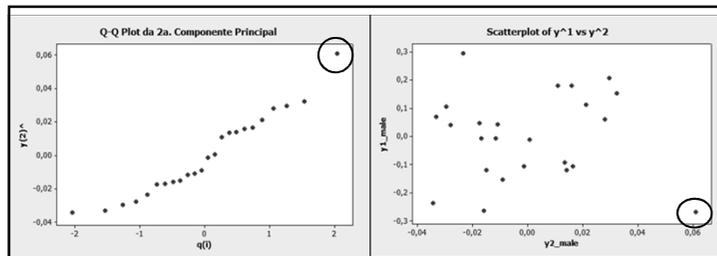
Let C30 = (C30-0,5)/24 # Cálculo percentagens

Name C31 "q(i)"
Invcdf c30 c31; # Cálculo quantis
Normal 0 1.

Name C32 "y(2)^"
Sort c25 c32 # Ordenação vetor de dados

Plot C32*C31; # Scatter plot
Title "Q-Q Plot da 2ª. Componente Principal";
Symbol.
```

Técnicas Multivariadas em Saúde - 2015



- Observação da 1ª. tartaruga é suspeita.
  - ✓ Checar registros ou verificar anomalias na tartaruga
- Excetuado esse dado o scatter plot aparenta estar razoavelmente elíptico
- Verificar os plots dos outros conjunto de componentes principais.

Técnicas Multivariadas em Saúde - 2015

### Propriedades Assintóticas

- Assuma que a amostra são observações aleatórias de população normal p-variada

✓ Autovalores desconhecidos são distintos e positivos

✓ Distribuição amostral autovalores

$$\sqrt{n}(\hat{\lambda} - \lambda) \stackrel{\text{as.}}{\sim} N_p(0, 2\Lambda^2) \quad \hat{\lambda}_i \stackrel{\text{as.}}{\sim} N\left(\lambda_i, 2\frac{\lambda_i^2}{n}\right)$$

✓ Distribuição amostral dos autovetores

$$\sqrt{n}(\hat{e}_i - e_i) \stackrel{\text{as.}}{\sim} N_p(0, E_i). \quad E_i = \lambda_i \sum_{k=1}^p \frac{\lambda_k}{\lambda_k - \lambda_i} e_k - e_k'$$

✓ Cada  $\hat{\lambda}_i$  é independente dos elementos de  $\hat{e}_i$  associados

Técnicas Multivariadas em Saúde - 2015

- Intervalo de confiança aproximado para os  $\lambda_i$  de amostras suficientemente grandes

$$\hat{\lambda}_i \stackrel{\text{as.}}{\sim} N\left(\lambda_i, 2\frac{\lambda_i^2}{n}\right) \quad P\left\{|\hat{\lambda}_i - \lambda_i| \leq z_{\alpha/2} \lambda_i \sqrt{\frac{2}{n}}\right\} = 1 - \alpha$$

$$\frac{\hat{\lambda}_i}{1 + z_{\alpha/2} \sqrt{\frac{2}{n}}} \leq \lambda_i \leq \frac{\hat{\lambda}_i}{1 - z_{\alpha/2} \sqrt{\frac{2}{n}}}$$

- Intervalos de confiança simultâneos de Bonferroni para  $m \lambda_i$ 's

√ Trocar  $z_{\alpha/2}$  por  $z_{\alpha/2m}$ .

Técnicas Multivariadas em Saúde - 2015

### Componentes Principais para Matrizes de Covariâncias com Estruturas Especiais

- Matriz diagonal:  $\Sigma_{p \times p} = \begin{bmatrix} \sigma_{11} & 0 & \dots & 0 \\ 0 & \sigma_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{pp} \end{bmatrix}$ .

√ j-ésimo autovetor:  $e'_j = [0, \dots, 0, 1, 0, \dots, 0]$

- 1 na j-ésima posição

√ j-ésima componente principal:  $Y_j = e'_j \mathbf{X} = X_j$

√ Não há ganho extraíndo as componentes principais

- A padronização não altera a situação

$$\mathbf{P} = \mathbf{I}$$

Técnicas Multivariadas em Saúde - 2015

### Exercício – Solo

- Análise de solo

√ 20 amostras

√ Variáveis:

- areia (%)
- sedimentos (%)
- argila (%)
- qte. material orgânico (%)
- acidez do solo (pH)

√ Banco de dados: *BD\_multivariada.xls/solo*

Técnicas Multivariadas em Saúde - 2015

- Matriz de covariâncias amostral (S)

Covariâncias: areia; sedimentos; argila; morganico; ph					
	areia	sedimentos	argila	morganico	ph
areia	138,32674				
sedimentos	-102,12274	79,73818			
argila	-36,20400	22,38455	13,81945		
morganico	-0,94221	1,52661	-0,58439	0,64345	
ph	-0,13579	0,11079	0,02500	0,03237	0,26263

- Autovalores de S

Eigenvalues				
223,841	8,218	0,472	0,258	0,000

√ S é singular pois  $\lambda_5 = 0$  ( $|S| = 0$ )

$$(X_1 + X_2 + X_3 = 100\%)$$

Técnicas Multivariadas em Saúde - 2015

• Componentes principais ( $p=5$ )

Principal Component Analysis: areia; sedimentos; argila; morganico; ph

Eigenanalysis of the Covariance Matrix

Eigenvalue	223,84	8,22	0,47	0,26	0,00
Proportion	0,962	0,035	0,002	0,001	0,000
Cumulative	0,962	0,997	0,999	1,000	1,000

Variable	PC1	PC2	PC3	PC4	PC5
areia	-0,785	0,223	-0,027	-0,004	-0,577
sedimentos	0,587	0,561	-0,086	-0,010	-0,577
argila	0,198	-0,788	0,113	0,014	-0,577
morganico	0,007	0,146	0,980	0,136	0,000
ph	0,001	0,002	0,137	-0,991	-0,000

- √  $y_5$  é constante para qualquer observação j  
 $y_5 = 0,577 (100)$
- √ Qualquer das três variáveis poderia ser eliminada

Técnicas Multivariadas em Saúde - 2015

- Eliminada  $X_1$  (areia)  
√ maior variância amostral  
tenderia dominar primeira componente
- Matriz de covariâncias amostral (S)

Covariâncias: sedimentos; argila; morganico; ph

	sedimentos	argila	morganico	ph
sedimentos	79,7382	22,3846	1,52661	0,110789
argila	22,3846	13,8194	-0,58439	0,025000
morganico	1,5266	-0,5844	0,64345	0,032368
ph	0,1108	0,0250	0,03237	0,262632

- Autovalores de S

Eigenvalues

86,6403	7,0936	0,4714	0,2584
---------	--------	--------	--------

Técnicas Multivariadas em Saúde - 2015

• Componentes principais ( $p = 4$  – eliminada  $X_1$ )

Principal Component Analysis: sedimentos; argila; morganico; ph

Eigenanalysis of the Covariance Matrix

Eigenvalue	86,640	7,094	0,471	0,258
Proportion	0,917	0,075	0,005	0,003
Cumulative	0,917	0,992	0,997	1,000

Variable	PC1	PC2	PC3	PC4
sedimentos	0,956	-0,288	0,059	0,006
argila	0,294	0,945	-0,142	-0,018
morganico	0,015	-0,154	-0,978	-0,136
ph	0,001	-0,002	-0,137	0,991

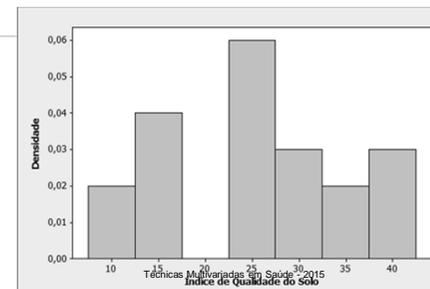
- √ Duas primeiras componentes explicam 99,2% da variância total
- 1ª. Componente: Índice de qualidade do solo em termos de % sedimentos e argila
  - sedimentos é a variável mais importante
- 2ª. Componente: Comparação entre % de sedimentos e % de argila
  - argila tem peso maior na componente
- 3ª. Componente: variável material orgânico

Técnicas Multivariadas em Saúde - 2015

- Scores da 1ª. Componente  
√ Índice de qualidade do solo em termos de sedimentos e areia  
(material orgânico tem pouco participação)

Descriptive Statistics: Scores y1

Variable	N	N*	Mean	SE Mean	StDev	Minimum	Q1	Median	Q3	Maximum
Scores Y1	20	0	24,86	2,08	9,31	11,57	14,77	24,29	32,43	39,70



Técnicas Multivariadas em Saúde - 2015

- Diferença de escala e unidades da variáveis  
 $\sqrt{\text{Recomendável padronização para análise de componentes}}$

Técnicas Multivariadas em Saúde - 2015

- Componentes principais (p=4) – Matriz de correlação

**Principal Component Analysis: sedimentos; argila; morganico; ph**

Eigenanalysis of the Correlation Matrix

Eigenvalue	1,6757	1,1461	0,9601	0,2181
Proportion	0,419	0,287	0,240	0,055
Cumulative	0,419	0,705	0,945	1,000

Variable	PC1	PC2	PC3	PC4
sedimentos	0,710	0,182	-0,147	-0,664
argila	0,702	-0,241	0,111	0,661
morganico	0,025	0,836	-0,423	0,349
ph	0,042	0,459	0,887	-0,026

Técnicas Multivariadas em Saúde - 2015

## Referências

Técnicas Multivariadas em Saúde - 2015

## Bibliografia Recomendada

- MANLY, B. J. F. *Métodos Estatísticos Multivariados: uma Introdução*. Bookman, 2008.
- JOHNSON, R. A.; WINCHERN, D. W. *Applied Multivariate Statistical Analysis*. Prentice Hall, 2007
- MINGOTI, D.C. *Análise de Dados através de Métodos de Estatística Multivariada*. Ed. UFMG, 2005.
- EVERITT, B.; HOTHORN, T. *An Introduction to Applied Multivariate Analysis with R*. Springer, 2011.

Técnicas Multivariadas em Saúde - 2015