

Análise Multivariada

Lupércio França Bessegato
Dep. Estatística/UFJF

Roteiro

1. Introdução
2. Representação de Dados Multivariados
3. Análise de Componentes Principais
4. Distribuições de Probabilidade Multivariadas
5. Análise Fatorial
6. Análise de Correlação Canônica
7. Análise de Conglomerados
8. Análise Discriminante
9. Referências

Análise Multivariada - 2016

2

Análise de Componentes Principais

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

Análise Multivariada - 2016

4

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

Análise Multivariada - 2016

5

- 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

Análise Multivariada - 2016

6

- 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

Análise Multivariada - 2016

7

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 Σ , cujos autovalores são $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$.

- Componentes principais de Σ :

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

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

Análise Multivariada - 2016

8

• Definição – 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:

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

√ 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 $(\mathbf{a}' \mathbf{X})$ é igual a λ_i .
- É obtida quando $\mathbf{a}_i = \mathbf{e}_i$.

• Definição 2 – Componente principal:

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

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

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

• Covariância entre duas componentes principais:

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

• Comentário:

√ Cada autovalor λ_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

- Variâncias total e generalizada de Σ :

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

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

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

Análise Multivariada - 2016

15

- 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

Análise Multivariada - 2016

16

- Aproximação de Σ :

√ 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 $p \times p$ correspondente apenas à informação da j-ésima componente principal

Análise Multivariada - 2016

17

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

Análise Multivariada - 2016

18

Estimação das Componentes Principais – Matriz de Covariâncias

- Em geral, Σ 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 : $\hat{e}_1, \hat{e}_2, \dots, \hat{e}_p$

- 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:

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

- 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

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

- Scores das componentes

√ Valor das componentes para cada elemento amostral

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

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*

- 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?

Análise Multivariada - 2016

29

| | 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

Análise Multivariada - 2016

30

- 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

Análise Multivariada - 2016

31

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

Análise Multivariada - 2016

32

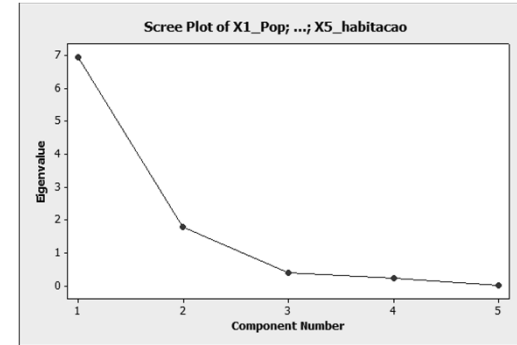
Scree Plot

- Gráfico λ_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

Análise Multivariada - 2016

34

Exemplo 8.3



Análise Multivariada - 2016

35

Exemplo 8.4

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

Análise Multivariada - 2016

36

Vetor de médias amostral (\bar{x})

```
Descriptive Statistics: log_comp_male; log_larg_male; log_esp_male
Variable      Mean
log_comp_male 4,7254
log_larg_male  4,8776
log_esp_male   3,7032
```

Matriz de covariâncias amostral (S)

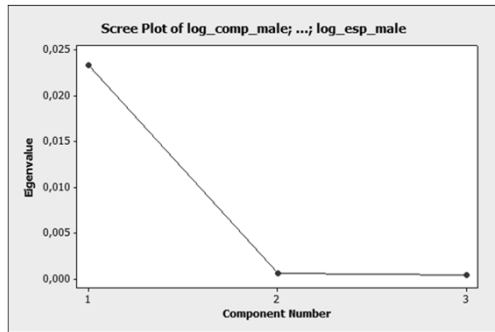
```
Covariances: log_comp_male; log_larg_male; log_esp_male
log_comp_male  log_comp_male  log_larg_male  log_esp_male
log_comp_male  0,01107200
log_larg_male  0,00801914  0,00641673
log_esp_male  0,00815965  0,00600527  0,00677276
```

- A variação amostral pode ser resumida por uma componente principal?

Análise Multivariada - 2016

37

• Scree Plot



√ Uma componente principal é claramente dominante

• 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

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} .

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

Comentários

- As componentes principais de Σ 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

Análise Multivariada - 2016

47

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

Análise Multivariada - 2016

48

- 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:

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

Análise Multivariada - 2016

49

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

Análise Multivariada - 2016

51

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 **R**: $\hat{e}_1, \hat{e}_2, \dots, \hat{e}_p$

- **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:

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

Exemplo 8.5

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

√ Período: Jan./75 a Dez./76

√ Ações:

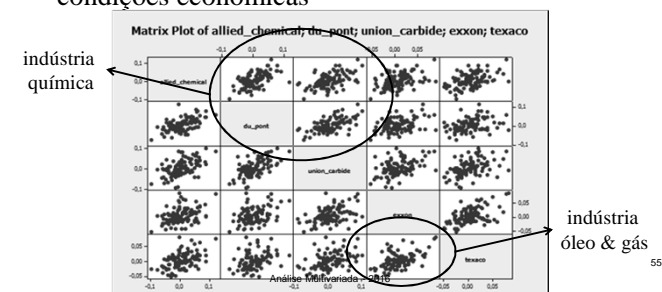
- Allied Chemical
- du Pont
- Union Carbide
- Exxon
- Texaco

√ Dados: BD_multivariada.xls/

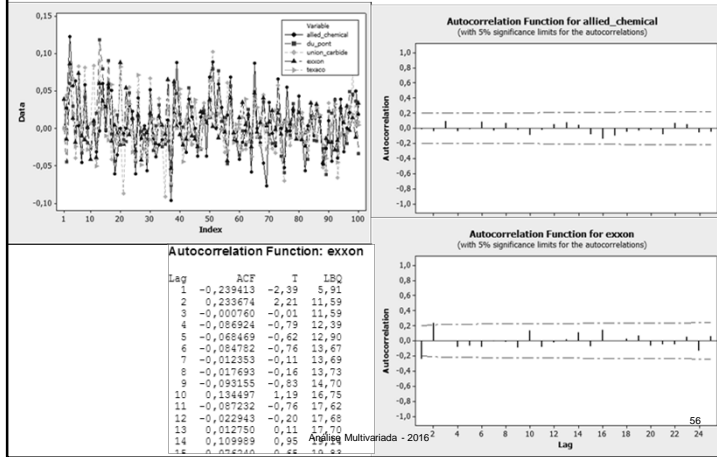
- Taxa de retorno: $\left(\text{taxa de retorno} \right) = \frac{\left(\text{\$ fechamento} \right)_{\text{sexta atual}} - \left(\text{\$ fechamento} \right)_{\text{sexta anterior}}}{\left(\text{\$ fechamento} \right)_{\text{sexta anterior}}}$

- 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 (R)

Correlations: allied_chemical; du_pont; union_carbide; Exxon; texaco

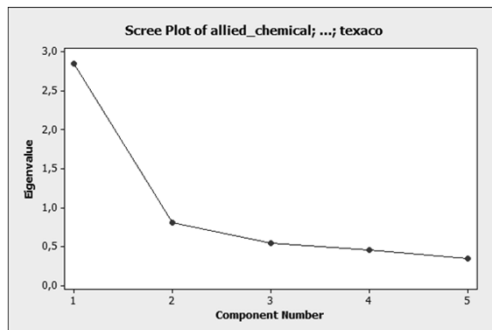
| | allied_chemical | du_pont | union_carbide | Exxon |
|---------------|-----------------|---------|---------------|-------|
| du_pont | 0,577 | | | |
| union_carbide | 0,509 | 0,598 | | |
| Exxon | 0,387 | 0,390 | 0,436 | |
| texaco | 0,462 | 0,322 | 0,426 | 0,524 |

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

Análise Multivariada - 2016

57

- Scree Plot



√ Aparentemente duas componentes principais resumem bem os dados

Análise Multivariada - 2016

58

- Componentes principais:

Principal Component Analysis: allied_chemi; du_pont; union_carbid; Exxon; texac

Eigenanalysis of the Correlation Matrix

| | | | | | |
|------------|--------|--------|--------|--------|--------|
| Eigenvalue | 2,8565 | 0,8091 | 0,3400 | 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

Análise Multivariada - 2016

59

Slide 57

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)

Análise Multivariada - 2016

60

• 2ª. componente principal:

$$\hat{y}_2 = 0,240z_1 + 0,509z_2 + 0,260z_3 - 0,526z_4 - 0,582z_5$$

√ Interpretação:

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

Análise Multivariada - 2016

61

• 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

Análise Multivariada - 2016

62

Variáveis Padronizadas – Regra Empírica

- Reter apenas as componentes cujas variâncias (λ_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

Análise Multivariada - 2016

63

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)

Análise Multivariada - 2016

73

Gráfico dos Componentes Principais

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

Análise Multivariada - 2016

76

- 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^a.s componentes são aproximadamente normais
- As últimas componentes principais podem ajudar a apontar observações suspeitas

Análise Multivariada - 2016

77

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.

Análise Multivariada - 2016

81

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:

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

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

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

- Comandos Minitab para Q-Q Plot

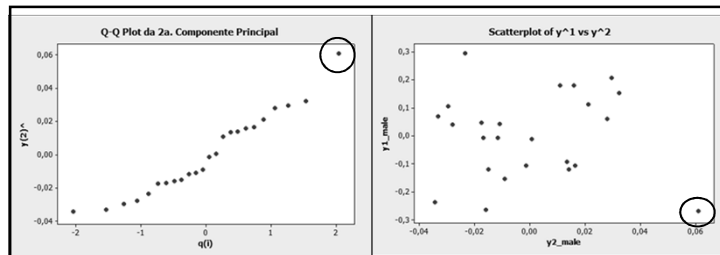
```
Name C30 "(j-1/2)/n"
Set C30
1(1 : 24 / 1)1
End.

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.
```



- 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.

Propriedades Assintóticas

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

✓ Autovalores desconhecidos são distintos e positivos

✓ Distribuição amostral autovalores

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

✓ Distribuição amostral dos autovetores

$$\sqrt{n}(\hat{\mathbf{e}}_i - \mathbf{e}_i) \stackrel{\text{as.}}{\sim} N_p(\mathbf{0}, \mathbf{E}_i). \quad \mathbf{E}_i = \lambda_i \sum_{\substack{k=1 \\ k \neq i}}^p \frac{\lambda_k}{(\lambda_k - \lambda_i)^2} \mathbf{e}_k \mathbf{e}_k'$$

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

- Intervalo de confiança aproximado para os λ_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 λ_i 's

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

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

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

Exemplo

- Estudo de poluição do ar em 41 cidades dos EUA

√ Ano: 1970

- Dados: *Usairpollution*{MVA}

- Variáveis:

√ SO2: conteúdo de dióxido de enxofre no ar, em $\mu\text{g}/\text{m}^3$.

√ Temp: temperatura média anual ($^{\circ}\text{F}$)

√ Indust: quantidade de empresas manufatureiras empregando pelo menos 20 empregados.

√ Pop: população (censo 1970), em milhares.

√ Vento: velocidade média anual de vento, em milhas/h

√ Precip: precipitação média anual, em polegadas

√ Dias: número médio anual de dias com precipitação

• Variáveis

- √ 'Resposta': SO2
- √ Ambientais: Vento, Precip, Dias, Temp
- √ Demográficas: Indust, Pop

• Comandos em R:

√ Carregamento dos dados:

```
> library(MVA)
> library(HSAUR2, ADGofTest)# carrega os pacotes
> data(USairpollution)# carrega o banco de dados
> poluicao <- USairpollution
> colunas<- c("SO2", "Temp", "Indust", "Pop",
+ "Vento", "Precip", "Dias")
> colnames(poluicao)<- colunas
> poluicao$negtemp <- poluicao$Temp * (-1)
> poluicao$Temp <- NULL
```

√ Comentários:

- Extrair componentes da matriz de correlações:
 - Variáveis estão em escalas muito distintas
- Ignora a variável SO2 (considera só as 'explicativas')
- Uso da temperatura negativa
 - As 6 variáveis têm valores altos, de maneira que representam um ambiente menos atrativo

√ Matriz de correlações:

```
> # Descrição do conjunto de dados
>
> poluicao.corr <- cor(poluicao[,-1])
> poluicao.corr
      Indust      Pop      Vento      Precip      Dias      negtemp
Indust  1.00000000  0.95526935  0.23799683 -0.03241688  0.13182930  0.19004216
Pop     0.95526935  1.00000000  0.21264375 -0.02611873  0.04208319  0.06267813
Vento   0.23799683  0.21264375  1.00000000 -0.01299438  0.16410559  0.34973963
Precip  -0.03241688 -0.02611873 -0.01299438  1.00000000  0.49609671 -0.38625342
Dias    0.13182930  0.04208319  0.16410559  0.49609671  1.00000000  0.43024212
negtemp 0.19004216  0.06267813  0.34973963 -0.38625342  0.43024212  1.00000000
```

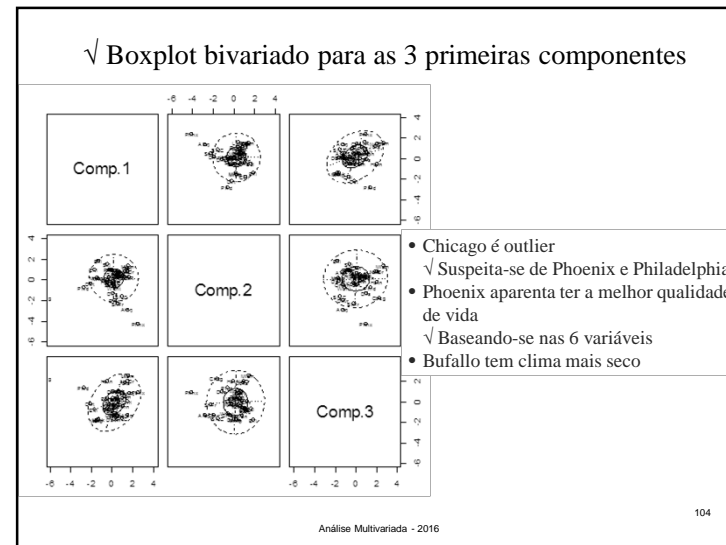
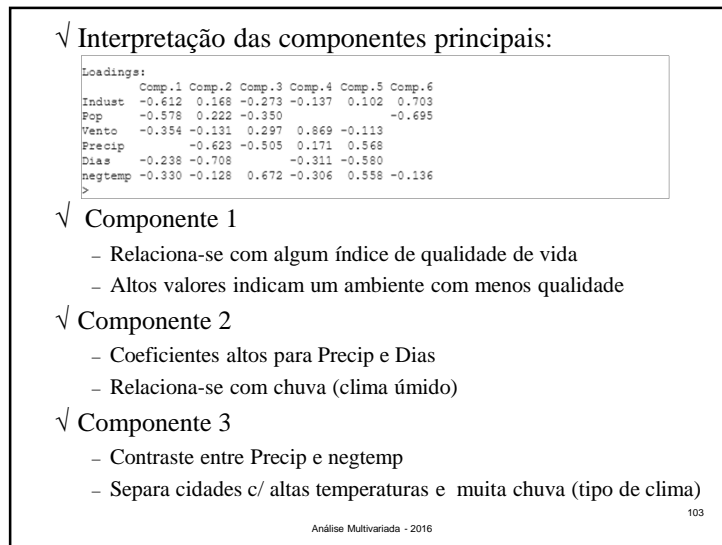
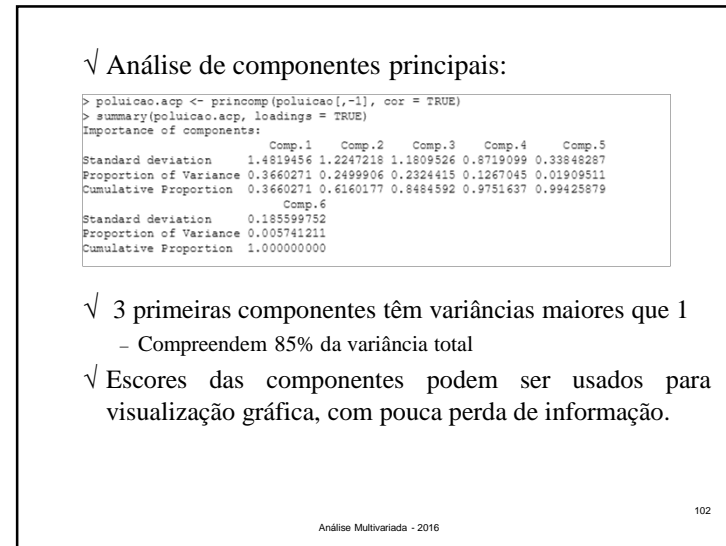
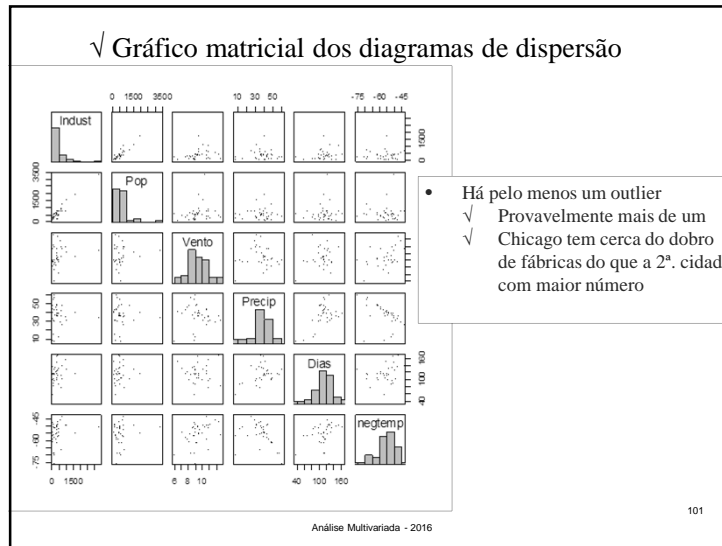
√ Comentário:

- Valores altos de correlação entre Indust e Pop

√ Matrix plot:

- Comandos em R:

```
> panel.hist <- function(x, ...){
+   usr <- par("usr"); on.exit(par(usr))
+   par(usr = c(usr[1:2], 0, 1.5) )
+   h <- hist(x, plot = FALSE)
+   breaks <- h$breaks; nb <- length(breaks)
+   y <- h$counts; y <- y/max(y)
+   rect(breaks[-nb], 0, breaks[-1], y, col = "grey", ...)
+ }
>
> pairs(poluicao[,-1], pch = ".", cex = 1.5)
> pairs(poluicao[,-1], diag.panel = panel.hist, pch = ".", cex = 1.5)
```



Questão Interessante

- Quais dentre as variáveis climáticas e ambientais são as melhores preditoras do grau de poluição do ar (concentração de SO₂)?

√ Esta questão é tratada com regressão linear múltipla

√ Potencial problema para aplicação dessa técnica:

- Alta correlação entre Indust e Pop

√ Solução:

- Retirar uma das variáveis

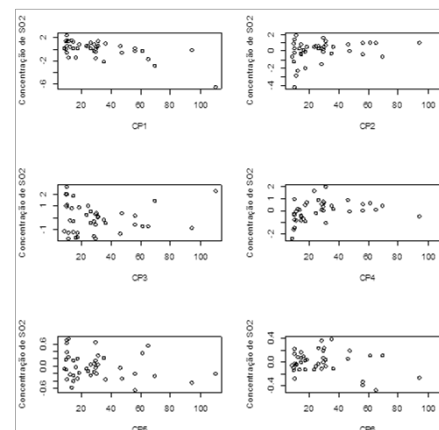
√ Alternativa:

- Fazer regressão dos níveis de SO₂ com as componentes principais derivadas das 6 variáveis originais
- Pode ser melhor regredir com todas as 6 componentes

Análise Multivariada - 2016

105

√ SO₂ dependendo das componentes principais



Análise Multivariada - 2016

106

√ Regressão com as 6 componentes principais:

```
> poluicao.reg <- lm(SO2 ~ poluicao.acp$scores,
+ data = poluicao)
> summary(poluicao.reg)

Call:
lm(formula = SO2 ~ poluicao.acp$scores, data = poluicao)

Residuals:
    Min       1Q   Median       3Q      Max
-23.004  -8.542  -0.991   5.758  48.758

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    30.049      2.286  13.146 6.91e-15 ***
poluicao.acp$scoresComp.1  -9.942      1.542  -6.446 2.28e-07 ***
poluicao.acp$scoresComp.2   2.240      1.866   1.200 0.23845
poluicao.acp$scoresComp.3   0.375      1.935   0.194 0.84752
poluicao.acp$scoresComp.4   8.549      2.622   3.261 0.00253 **
poluicao.acp$scoresComp.5 -15.176      6.753  -2.247 0.03122 *
poluicao.acp$scoresComp.6 -39.271     12.316  -3.189 0.00306 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 14.64 on 34 degrees of freedom
Multiple R-squared:  0.6695, Adjusted R-squared:  0.6112
F-statistic: 11.48 on 6 and 34 DF, p-value: 5.419e-07
```

√ Escores da 1ª. componente predizem mais a resposta

√ Componentes com menor variância não têm necessariamente as menores correlações com a resposta

Análise Multivariada - 2016

107

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*

Análise Multivariada - 2016

110

• 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 $\lambda_5 = 0$ ($|S| = 0$)
 $(X_1 + X_2 + X_3 = 100\%)$

• 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 | -0,794 | 0,113 | 0,014 | -0,577 |
| morganico | 0,007 | 0,146 | 0,380 | 0,136 | 0,000 |
| ph | 0,001 | 0,002 | 0,137 | -0,991 | -0,000 |

√ y5 é constante para qualquer observação j
 $y_5 = 0,577$ (100)

√ Qualquer das três variáveis poderia ser eliminada

• 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 |
|---------|--------|--------|--------|

• Componentes principais (p = 4 – eliminada X_1)

Principal Component Analysis: sedimentos; argila; morganico; ph

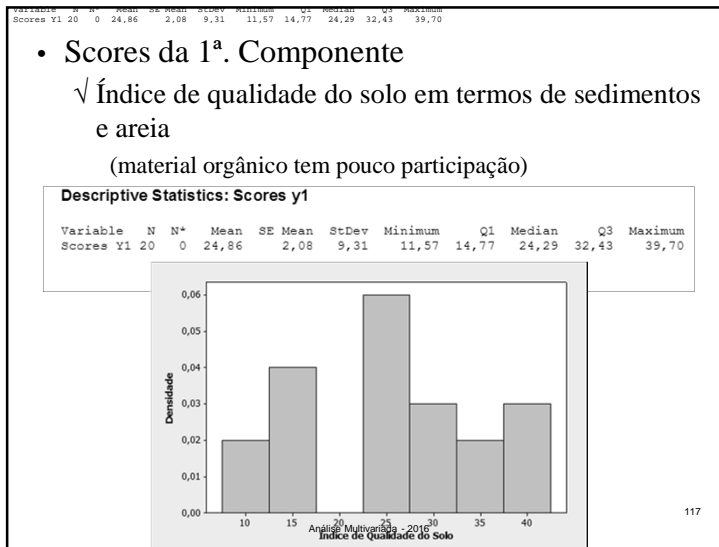
Eigenanalysis of the Covariance Matrix

| | | | | |
|------------|--------|-------|-------|-------|
| 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 | -0,184 | -0,379 | -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



- Diferença de escala e unidades da variáveis
 - √ Recomendável padronização para análise de componentes

118

- 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 | 0,023 | 0,836 | -0,423 | 0,349 |
| ph | 0,042 | 0,459 | 0,887 | -0,026 |

119

Referências

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.