

Introdução à Análise Multivariada usando R

Marcio Nicolau

2016-10-23

Contents

1	Pré-requisitos	5
2	Introdução	7
3	Análise de Agrupamentos	9
3.1	Exemplo: Força de trabalho agrícola na UE (1993)	9
3.2	Processo de Agrupamento	12
3.3	Métodos de Agrupamento	12
3.4	Distâncias para cálculo de agrupamento $\{\#\#\text{AAdist}\}$	12
4	Análise de Componentes Principais	15
5	Análise Fatorial	17
5.1	Example one	17
5.2	Example two	17

Chapter 1

Pré-requisitos

Antes de iniciar o curso, faz-se necessário instalar as ferramentas e pacotes para uso.

Será utilizado o software *R version 3.3.1 (2016-06-21)* disponível em CRAN e o *RStudio v1.0.44 Preview* disponível em RStudio.

A carga horária total será de 20 horas (3 dias) nos quais os seguintes tópicos de Análise Multivariada serão abordados.

- Análise de Agrupamentos (CA)
- Análise de Componentes Principais (PCA)
- Análise Fatorial (FA)

Alguns exemplos serão realizados com dados disponíveis na literatura científica e, em certos momentos, serão utilizados dados dos próprios participantes.

Espera-se que ao final do curso o aluno seja capaz de entender o uso de cada técnica e aplicá-la de forma correta em seu campo/área de pesquisa.

Cabe lembrar que este é um curso introdutório e que de forma alguma os conteúdos e aplicações serão apresentadas de forma exaustiva.

Marcio Nicolau

Estatístico / Embrapa Trigo

Palmas/TO, Outubro de 2016

Chapter 2

Introdução

As técnicas de Análise Multivariadas oferecem aplicações em diversas áreas do conhecimento no desenvolvimento científico.

Geralmente são utilizadas em fase exploratória de análise de dados, onde se busca entender melhor as relações entre as variáveis (medições físicas ou observações) de certo evento sob estudo ou de interesse científico.

Há também aplicações em conjunto com outros métodos da estatística onde é possível obter validações ou testes de carácter conclusivo.

Durante este curso e, certamente limitados pelo tempo, serão abordados somente as técnicas de caractere exploratório com a finalidade de melhor explicar as relações intrínsecas entre os dados, reduzir a dimensão, entender fontes de variabilidade, criar grupos homogêneos de indivíduos/espécies.

Pode-se dizer que a Análise Fatorial (FA), a Análise de Componentes Principais (PCA) e a Análise de Cluster (CA) são processo que tem por objetivo reduzir a complexidade dos dados observados, bem como entender o modelo estrutural presente nos dados.

No caso do FA, o objetivo é o de identificar construções poucos constructos para explicar os dados observados. No caso de PCA, pode não ser simples redução de dimensão, mas a interpretação dos componentes.

Por fim, a Análise de Cluster (CA) pode também ser usada para criar grupos de variáveis com interesse de reduzir a complexidade dos dados por meio da formação de grupos menores e homogêneos.

Tecnicamente, o problema de redução de dados pode ser resolvido como uma decomposição do valor singular (SVD) da matriz original, embora a solução mais típica seja o uso de PCA nas matrizes de covariância e/ou correlação.

Chapter 3

Análise de Agrupamentos

Nesta seção serão utilizados as seguinte bibliotecas do R.

```
libs <- c('cluster', 'psych')
sapply(libs, require, character.only = TRUE)
```

```
## Loading required package: cluster
```

```
## Loading required package: psych
```

```
## cluster  psych
##      TRUE    TRUE
```

3.1 Exemplo: Força de trabalho agrícola na UE (1993)

Estes conjunto registra os dados da produção per capita e o percentual da população que trabalha na agricultura em cada país da UE em 1993.

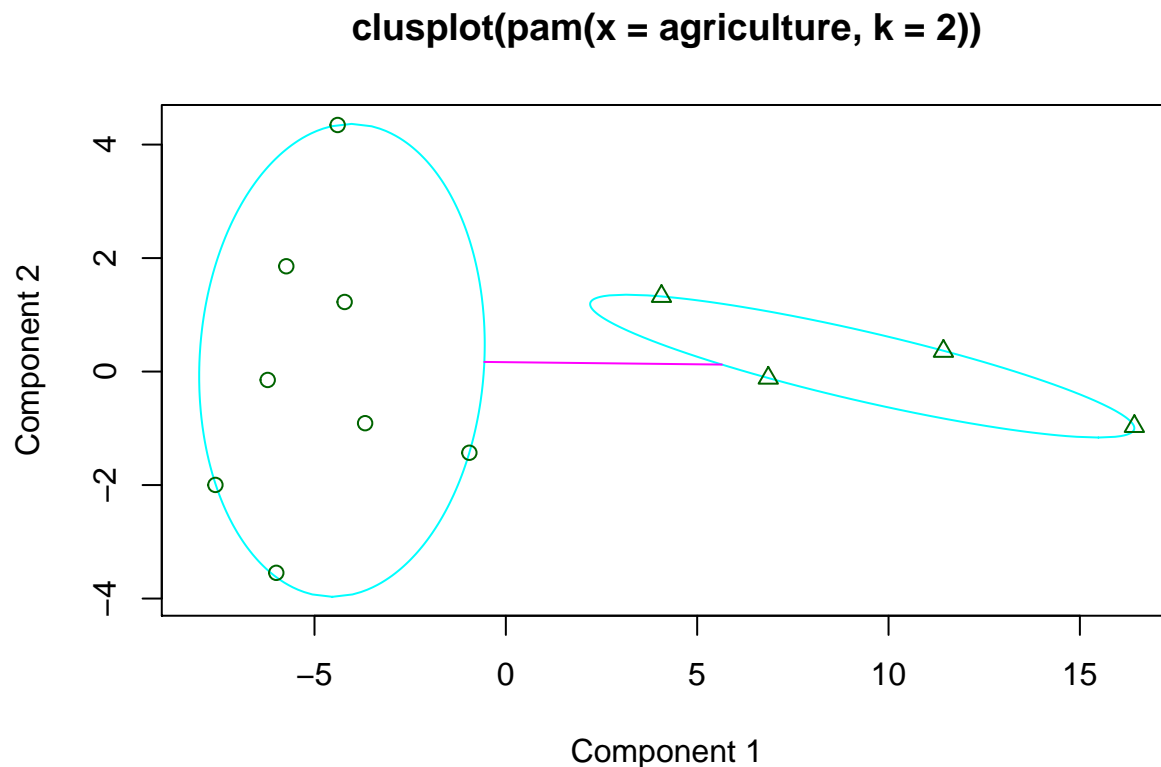
```
data(agriculture)
```

```
## Calcula matriz de dissimilaridade usando distância euclidiana
## e sem padronização das variáveis
daisy(agriculture, metric = "euclidean", stand = FALSE)
```

```
## Dissimilarities :
##           B           DK           D           GR           E           F           IRL
## DK    5.408327
## D     2.061553  3.405877
## GR    22.339651 22.570113 22.661200
## E      9.818350 11.182576 10.394710 12.567418
## F      3.448188  3.512834  2.657066 20.100995  8.060397
## IRL   12.747549 13.306014 13.080138  9.604166  3.140064 10.564563
## I      5.803447  5.470832  5.423099 17.383325  5.727128  2.773085  7.920859
## L      4.275512  2.220360  2.300000 24.035391 12.121056  4.060788 14.569145
## NL      1.649242  5.096077  2.435159 20.752349  8.280097  2.202272 11.150785
```

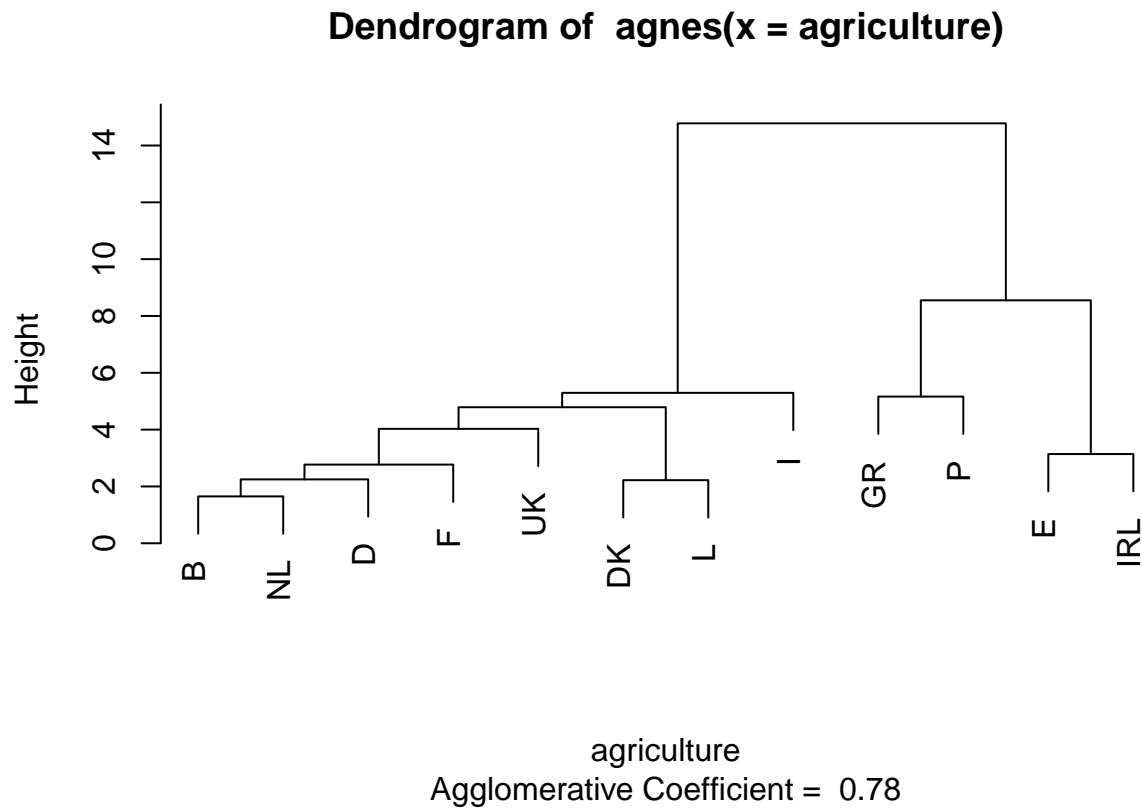
```
## P    17.236299 17.864490 17.664088  5.162364  7.430343 15.164432  4.601087
## UK    2.828427  8.052950  4.850773 21.485344  8.984431  5.303772 12.103718
##          I          L          NL          P
## DK
## D
## GR
## E
## F
## IRL
## I
## L    6.660330
## NL    4.204759  4.669047
## P    12.515990 19.168985 15.670673
## UK    6.723095  7.102112  3.124100 16.323296
##
## Metric :  euclidean
## Number of objects : 12
```

```
## Usa método de particionamento pelo meióide
## Partitioning Around Medoids (PAM)
plot(pam(agriculture, 2), which.plots = 1)
```



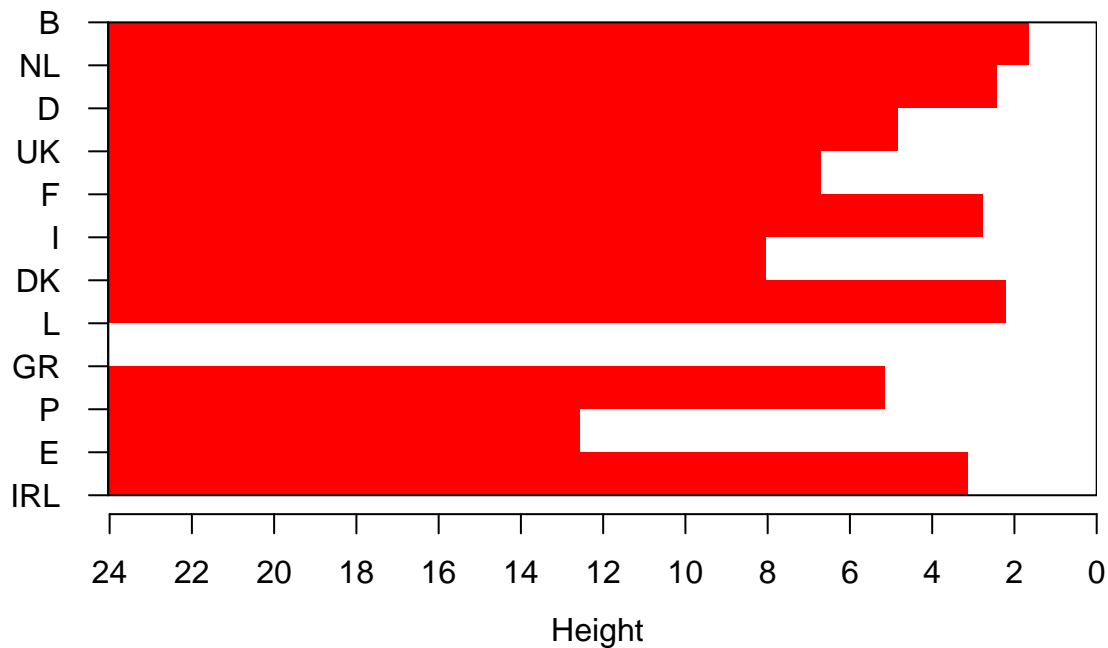
These two components explain 100 % of the point variability.

```
## Gráfico dendrograma usando método aglomeração mais próximo
## agnes
plot(agnes(agriculture), which.plots = 2)
```



```
## Plot dissimilaridade usando método divisivo  
## diana  
plot(diana(agriculture), which.plots = 1)
```

Banner of `diana(x = agriculture)`



Divisive Coefficient = 0.87

3.2 Processo de Agrupamento

Um agrupamento pode ser construído de duas formas:

- hierarquia: funções *agnes*, *diana*, *mona* e *hclust*;
- particionamento: funções *pam*, *clara*, *fanny* e *kmeans*

3.3 Métodos de Agrupamento

Um agrupamento pode gerar os grupos utilizando algum dos métodos a seguir (mais comuns):

- média: *average* ou UPGMA
- simple: *single*
- completa: *complete*
- Ward: *ward*
- média ponderada: *weighted* ou WPGMA

3.4 Distâncias para cálculo de agrupamento {##AAdist}

Para se calcular a distância entre os componentes, pode-se utilizar as funções a seguir (mais comuns):

- euclidiana: *euclidean*, raiz da soma dos quadrados das diferenças

- mahalanobis
- Manhattan: *manhattan*, soma das diferenças média absoluta
- Maximum:
- Canberra
- Binary
- Minkowski

Chapter 4

Análise de Componentes Principais

É uma alternativa à Análise Fatorial (FA), apesar dos objetivos serem semelhantes (PCA e FA), na PCA se busca obter o modelo descritivo dos dados enquanto na FA se busca o modelo estrutural.

Outro destaque importante é que a matriz/vetor de cargas “*loadings*” possuem valores equivalentes, na FA estes são menores. Isto ocorre porque na PCA é ajustado um modelo para a variância completa da matriz de correlação das variáveis e na FA o processo é realizado somente para a variância comum.

Chapter 5

Análise Fatorial

Some *significant* applications are demonstrated in this chapter.

5.1 Example one

5.2 Example two

Bibliography