

Clusterização Hierárquica

EST0133 - Introdução à Modelagem de Big Data

Marcus Nunes

<https://introbigdata.org/>

<https://marcusnunes.me/>

Universidade Federal do Rio Grande do Norte

Introdução

Introdução

- Outra maneira de clusterizar dados
- Enquanto o k-means se utiliza de somas de quadrados e distância euclidiana, a clusterização hierárquica se baseia em outras distâncias e similaridades
- É outro algoritmo que também pode ser baseado no vizinho mais próximo
- Busca coesão interna (intra grupo) e isolamento externo (entre grupos)

Distância

Distância

- Distância é um conceito bem definido em matemática
- Uma distância em um espaço métrico X é qualquer função

$$d : X \times X \rightarrow [0, \infty)$$

que satisfaça as quatro condições abaixo para todo $\mathbf{x}, \mathbf{y}, \mathbf{z} \in X$:

- i) $d(\mathbf{x}, \mathbf{y}) \geq 0$
- ii) $d(\mathbf{x}, \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} = \mathbf{y}$
- iii) $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$
- iv) $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$

Distância

- Distância discreta: se $\mathbf{x} = \mathbf{y}$, então $d(\mathbf{x}, \mathbf{y}) = 0$. Caso contrário, $d(\mathbf{x}, \mathbf{y}) = 1$.
- Distância euclidiana: sejam $X = \mathbb{R}^k$ e sejam $\mathbf{x}, \mathbf{y} \in \mathbb{R}^k$ tais que $\mathbf{x} = (x_1, \dots, x_n)$ e $\mathbf{y} = (y_1, \dots, y_n)$. Assim,

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2} = \left[\sum_{i=1}^n (x_i - y_i)^2 \right]^{\frac{1}{2}}$$

Distância

- Distância de Manhattan:

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n |x_i - y_i|$$

- Distância de Minkowski (ou norma- p):

$$d(\mathbf{x}, \mathbf{y}) = \left[\sum_{i=1}^n |x_i - y_i|^p \right]^{\frac{1}{p}}$$

- Como Escolher a Distância a ser Usada?
- A distância euclidiana funciona muito bem em casos lineares, mas é sensível a outliers
- A distância Manhattan funciona melhor em casos com muitas observações aberrantes ou num espaço com muitas dimensões
- Outras distâncias podem ser usadas em outros casos
- Essa escolha é arbitrária

Ligação

- Outro critério a ser utilizado na clusterização hierárquica é o método de ligação
- Assim como a distância escolhida, o método de ligação influenciará na clusterização final
- Também é uma escolha arbitrária

Sejam $\mathbf{x} \in X$ e $\mathbf{y} \in Y$ vetores nos conjuntos X e Y e $d(\mathbf{x}, \mathbf{y})$ alguma distância entre estes vetores

- Simples: mínima distância ou vizinho mais próximo - $\min\{d(\mathbf{x}, \mathbf{y})\}$
- Completa: máxima distância ou vizinho mais distante - $\max\{d(\mathbf{x}, \mathbf{y})\}$
- Ward: mínima variância - $\min \|\mathbf{x} - \mathbf{y}\|^2$
- E muitos outros

Algoritmo

1. Para a clusterização aglomerativa, comece com n nós de singletons (isto é, de um sujeito por nó)
2. Calcule a similaridade dois a dois entre todos os pares de nós
3. Junte os dois nós mais similares entre si em um nó
4. Repita os passos 2. e 3. até obter o número desejado de conglomerados

- O pesquisador deve definir que tipo de medida de similaridade utilizar para esta análise
- É preciso decidir o melhor tipo de medida de distância/similaridade entre os sujeitos (distância euclidiana, correlação etc.) e método de ligação
- Em geral, os dados são apresentados em um dendrograma

Número de Clusters

Número de Clusters

- Assim como no caso do k-means, precisamos determinar o número correto de clusters
- É uma tarefa difícil, que é melhor executada com conhecimentos *a priori* sobre os dados
- Vamos ver como utilizar a estatística gap para escolher o número de clusters

Número de Clusters

- Seja W_k a estatística gap para k clusters
- W_k parte do pressuposto de que, se as unidades amostrais não formassem grupos, ela seria aproximadamente uniforme para qualquer número de grupos
- Se existem grupos nos dados, então W_k forma um cotovelo no número ótimo de grupos

Número de Clusters

- W_k é dada por

$$W_k = \sum_{r=1}^k \frac{1}{2n_r} D_r,$$

em que $D_r = \sum_{i,j \in C_r} d_{ij}$ e

- k é o número de clusters
- n_r é o número de observações em cada cluster C_r
- d_{ij} é a distância entre as observações i e j

Exemplo

Exemplo

- Utilizaremos os dados sobre prisões nos Estados Unidos
- Esse conjunto de dados traz as seguintes taxas:
 - **Murder**: prisões por assassinato (a cada 100.000 habitantes)
 - **Assault**: prisões por assalto (a cada 100.000 habitantes)
 - **UrbanPop**: percentual de população urbana
 - **Rape**: prisões por estupro (a cada 100.000 habitantes)
- Os dados são em nível estadual, coletados em 1973
- Veremos como combinar distância, métodos de ligação e número de clusters nas nossas análises

Exemplo

```
> library(tidyverse)
> library(ggdendro)
> library(factoextra)
>
> prisoes <- scale(USArrests)
>
> prisoes.euclidiana <- dist(prisoes, method = "euclidean")
>
> prisoes.manhattan <- dist(prisoes, method = "manhattan")
```

Exemplo

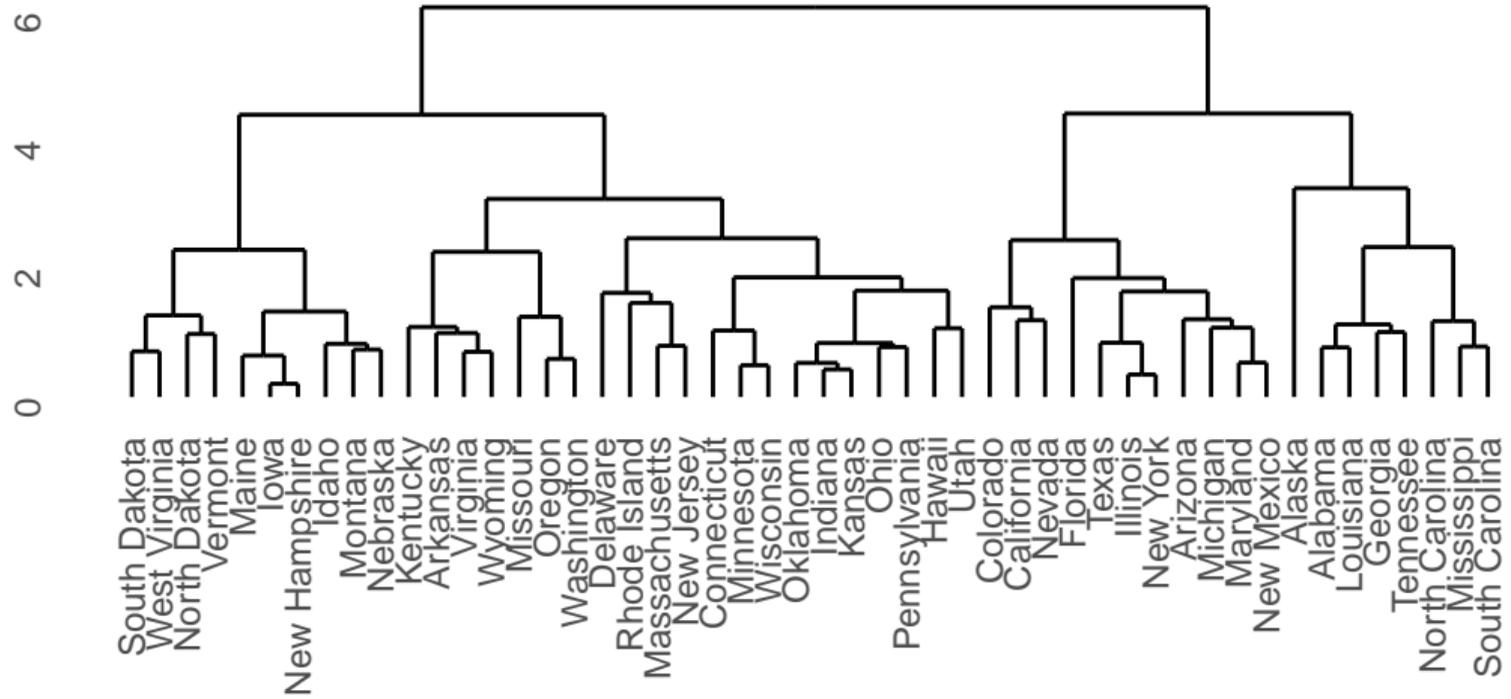
```
> ggdendrogram(hclust(priso.es.euclidiana, method = "single")) +  
+   labs(title = "Dist = Euclidiana, Link = Simples")  
>  
> ggdendrogram(hclust(priso.es.euclidiana, method = "complete")) +  
+   labs(title = "Dist = Euclidiana, Link = Completo")  
>  
> ggdendrogram(hclust(priso.es.euclidiana, method = "ward.D")) +  
+   labs(title = "Dist = Euclidiana, Link = Ward")
```

Exemplo

```
> ggdendrogram(hclust(priso.es.manhattan, method = "single")) +  
+   labs(title = "Dist = Manhattan, Link = Simples")  
>  
> ggdendrogram(hclust(priso.es.manhattan, method = "complete")) +  
+   labs(title = "Dist = Manhattan, Link = Completo")  
>  
> ggdendrogram(hclust(priso.es.manhattan, method = "ward.D")) +  
+   labs(title = "Dist = Manhattan, Link = Ward")
```

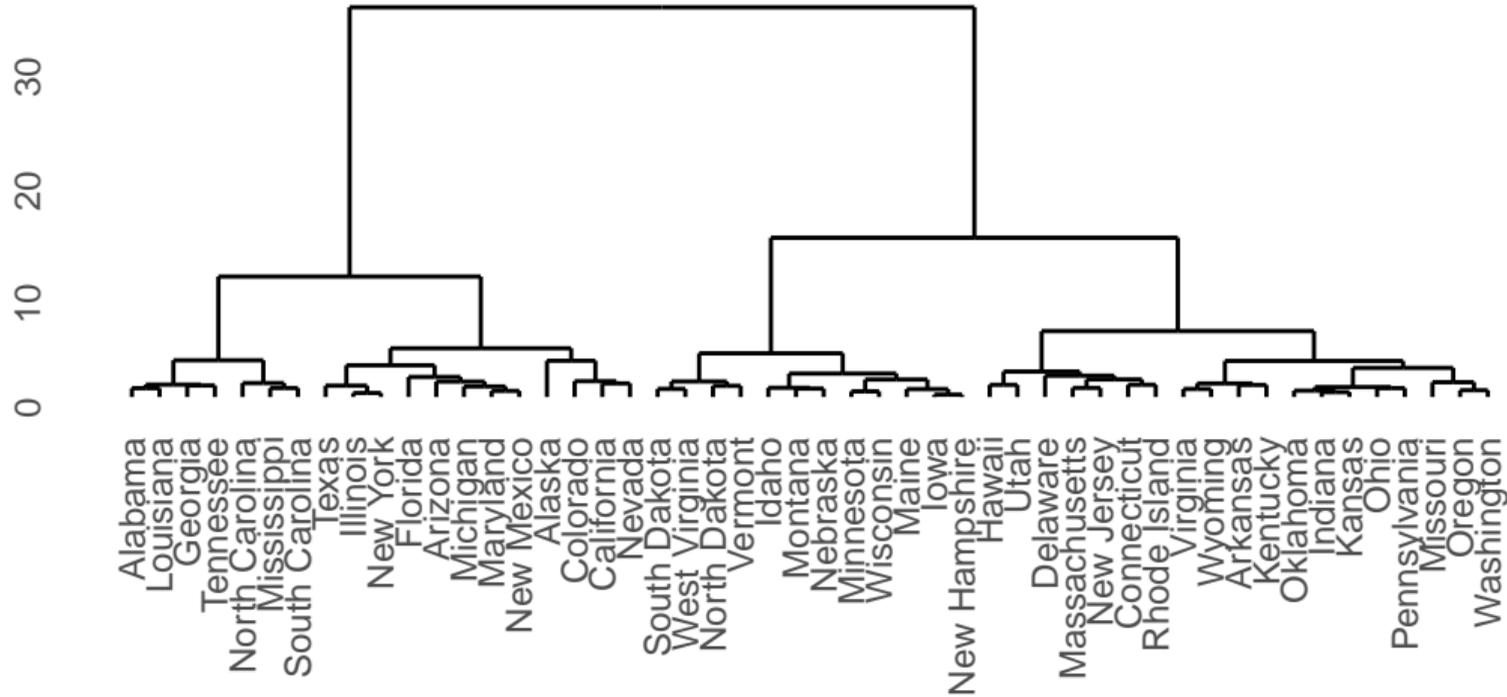

Exemplo

Dist = Euclidiana, Link = Completo



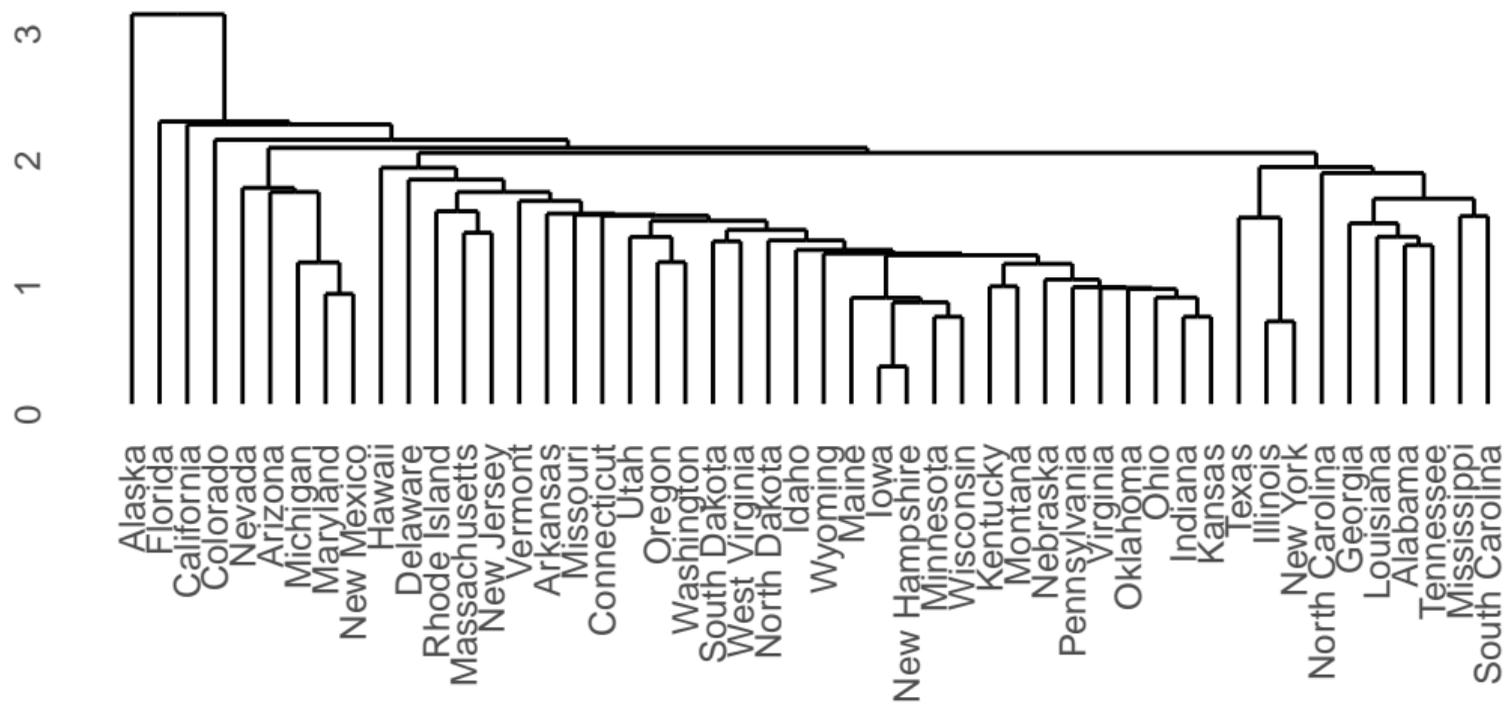
Exemplo

Dist = Euclidian, Link = Ward



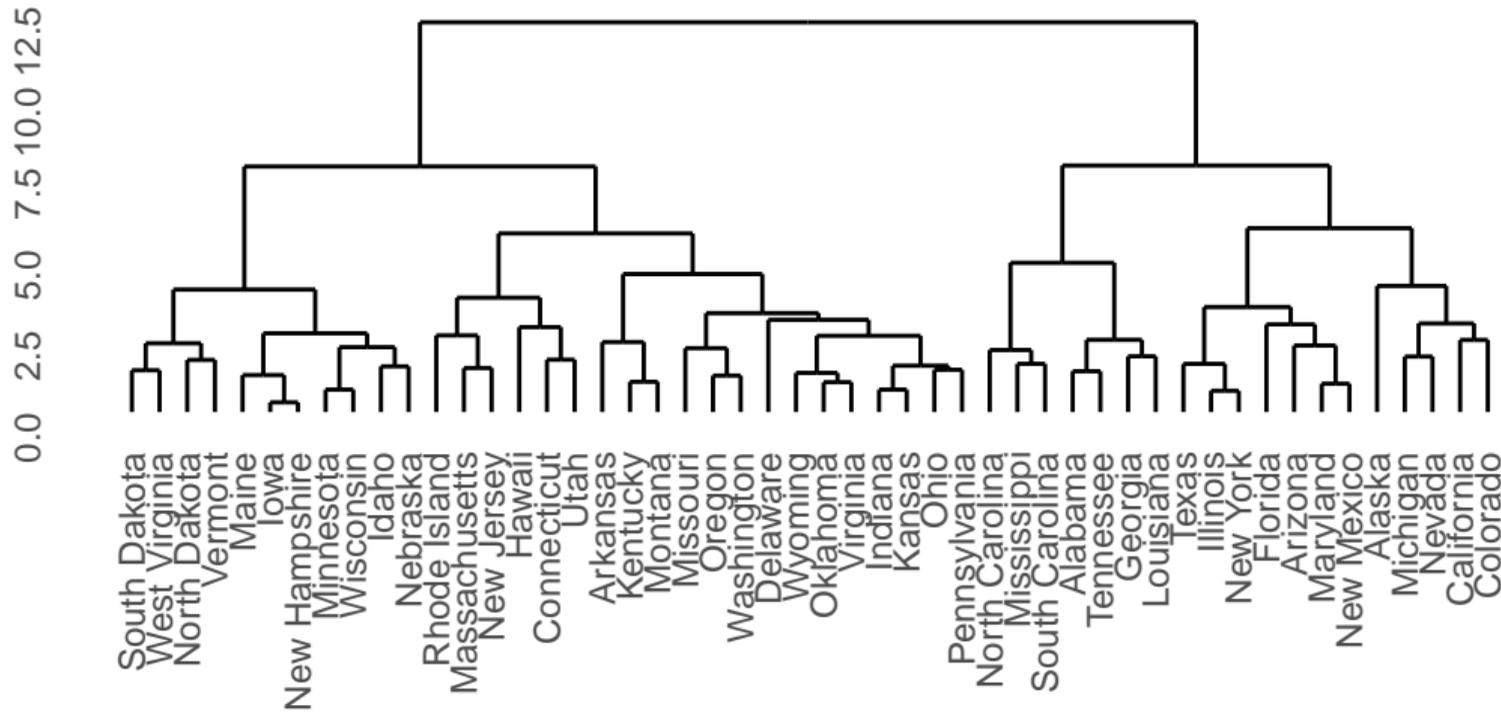
Exemplo

Dist = Manhattan, Link = Simples



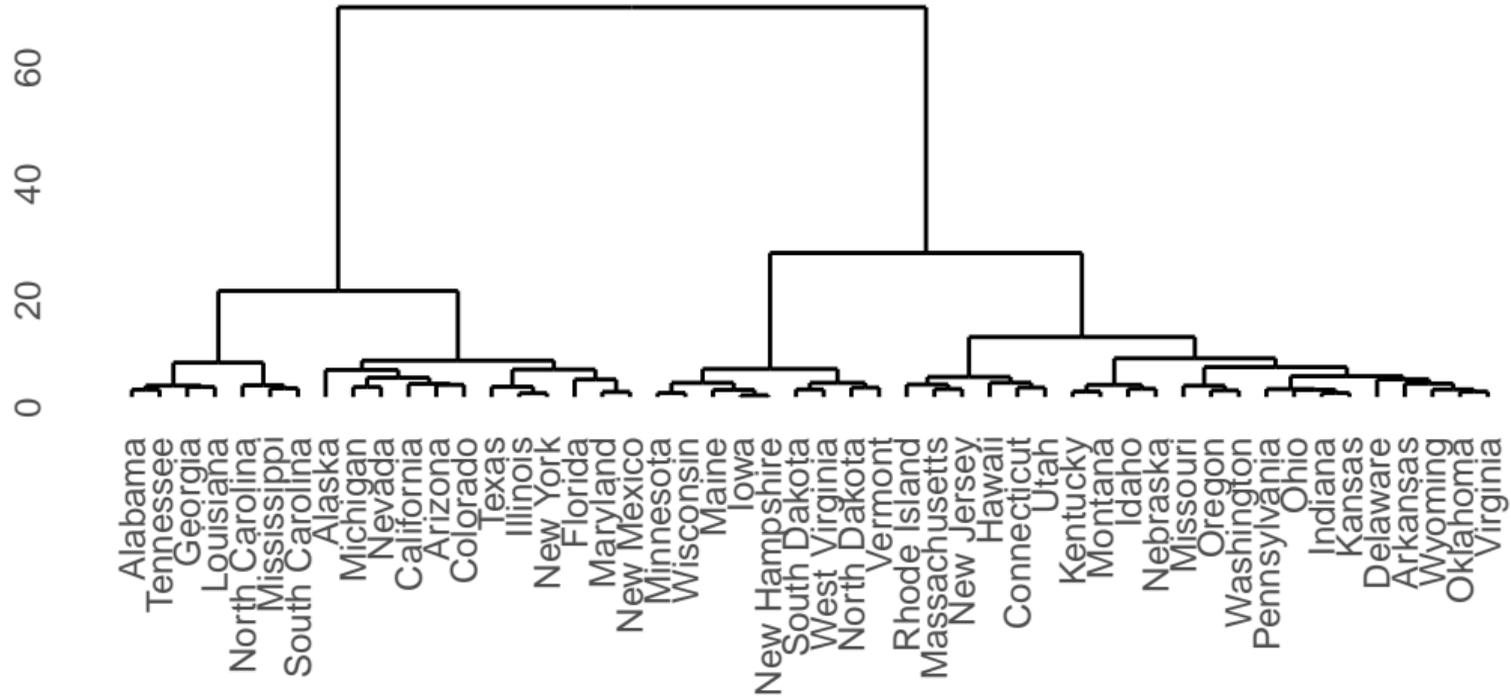
Exemplo

Dist = Manhattan, Link = Completo



Exemplo

Dist = Manhattan, Link = Ward



Exemplo

```
> euclidiana_single <- function(x, k) {  
+   list(cluster = cutree(hclust(dist(x),  
+                           method = "single"),  
+                           k=k))  
+ }  
>  
> fviz_nbclust(prisoos,  
+               FUNcluster = euclidiana_single,  
+               method = "gap_stat",  
+               nboot = 100)
```

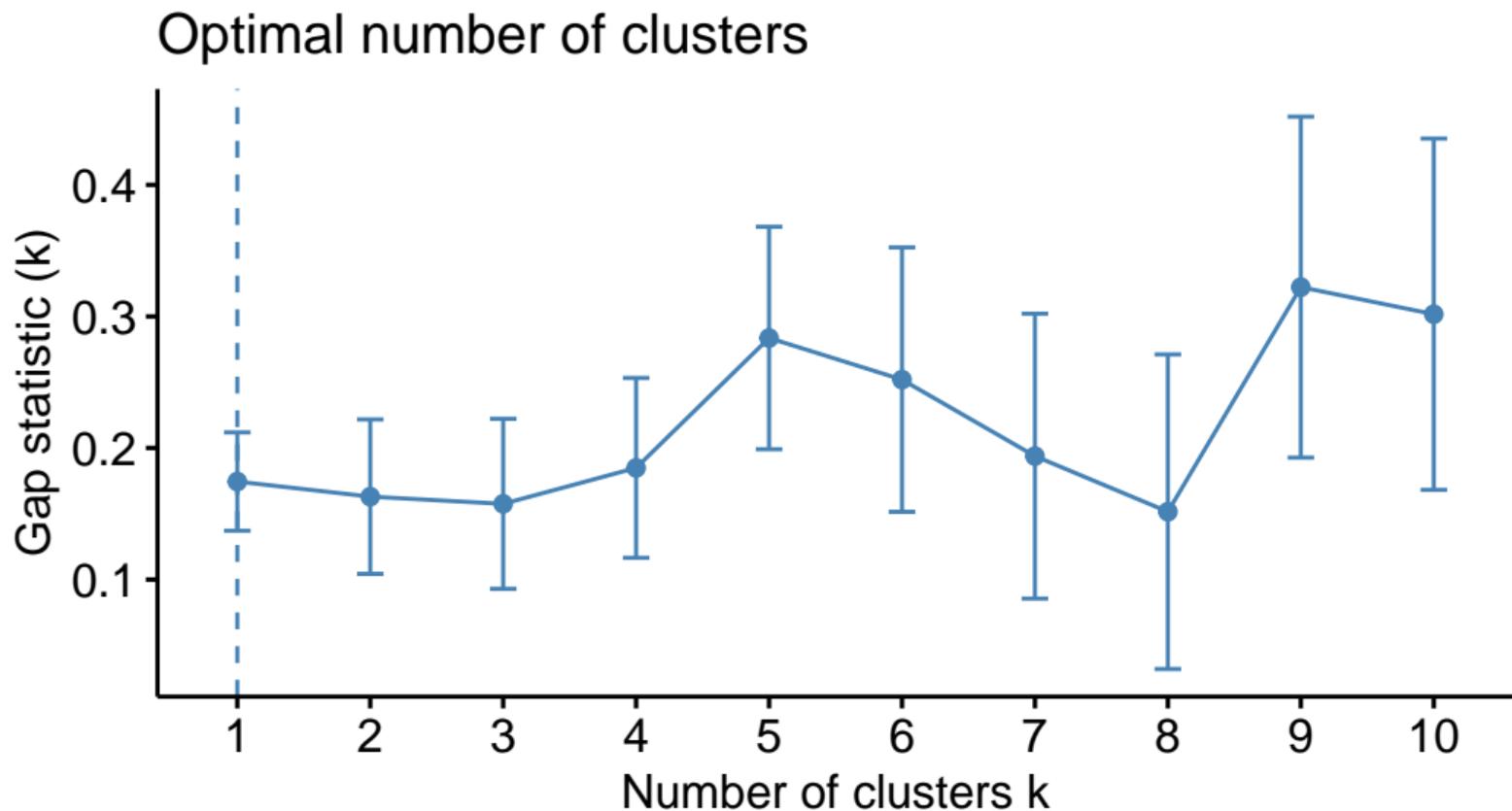
Exemplo

```
> euclidiana_complete <- function(x, k) {  
+   list(cluster = cutree(hclust(dist(x),  
+                           method = "complete"),  
+                           k=k))  
+ }  
>  
> fviz_nbclust(prisoos,  
+               FUNcluster = euclidiana_complete,  
+               method = "gap_stat",  
+               nboot = 100)
```

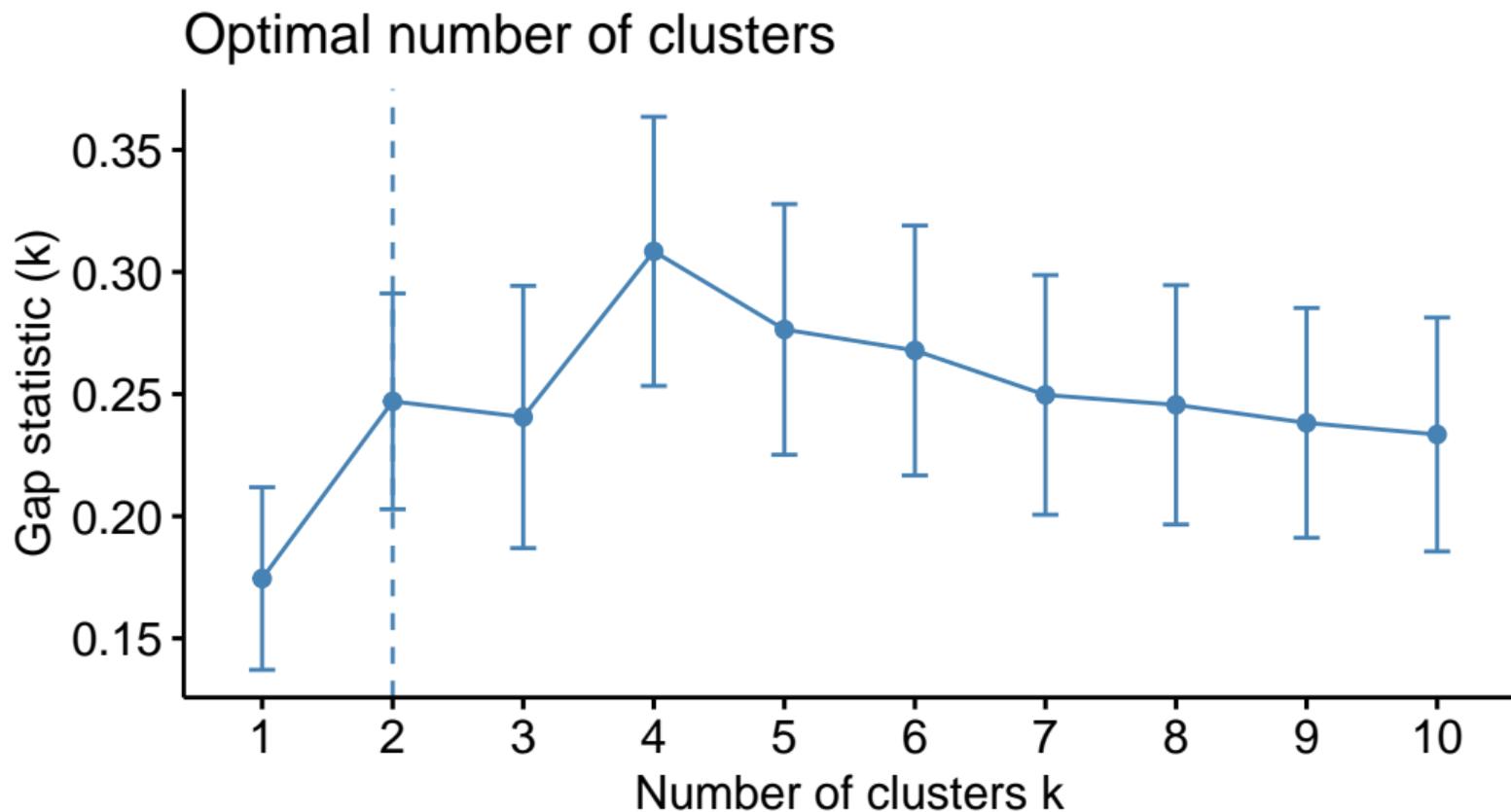
Exemplo

```
> euclidiana_ward <- function(x, k) {  
+   list(cluster = cutree(hclust(dist(x),  
+                           method = "ward.D"),  
+                           k=k))  
+ }  
>  
> fviz_nbclust(prisoos,  
+               FUNcluster = euclidiana_ward,  
+               method = "gap_stat",  
+               nboot = 100)
```

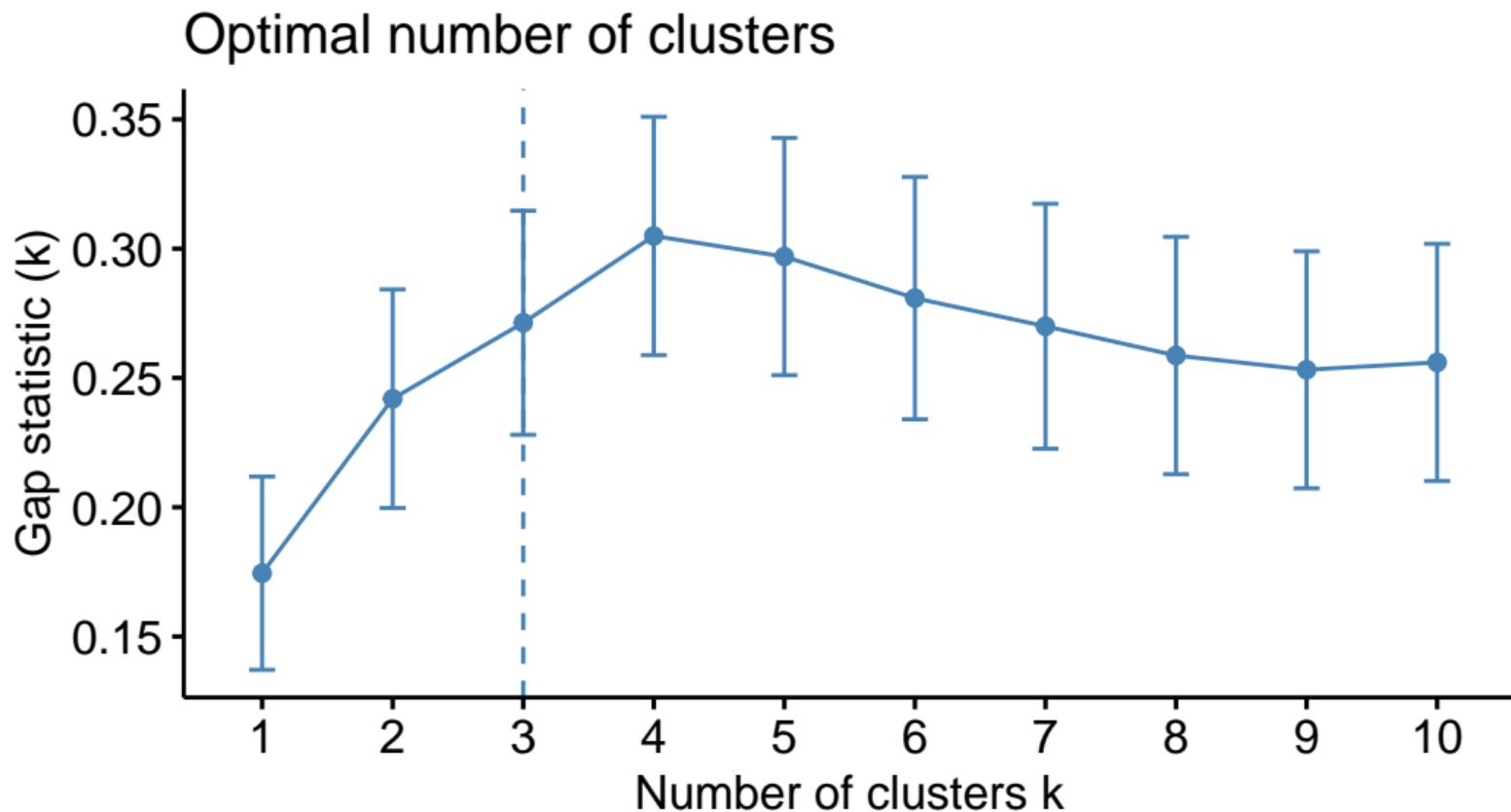
Exemplo



Exemplo



Exemplo



Exemplo

- O número correto de clusters vai ser definido de acordo com o interesse ou conhecimento do pesquisador
- A distância a ser utilizada depende dos dados
- Esta é outra ferramenta de análise exploratória

Exercícios

Exercícios

1. O arquivo `AlimentacaoReinoUnido.txt` mostra o consumo de diversos alimentos no Reino Unido em 1997. Importe este conjunto de dados para o R.
2. Faça a clusterização hierárquica deste conjunto de dados utilizando a distância euclidiana.
3. Faça a clusterização hierárquica deste conjunto de dados utilizando a distância de Manhattan.
4. Há algum padrão interessante nos dados? Qual ou quais?
5. Refaça a análise para o conjunto de dados com a transposta da matriz original de dados. O que mudou? O que é possível perceber com esta nova análise?

6. Repita os exercícios do slide anterior para o conjunto de dados `heptatlo`

Clusterização Hierárquica

EST0133 - Introdução à Modelagem de Big Data

Marcus Nunes

<https://introbigdata.org/>

<https://marcusnunes.me/>

Universidade Federal do Rio Grande do Norte