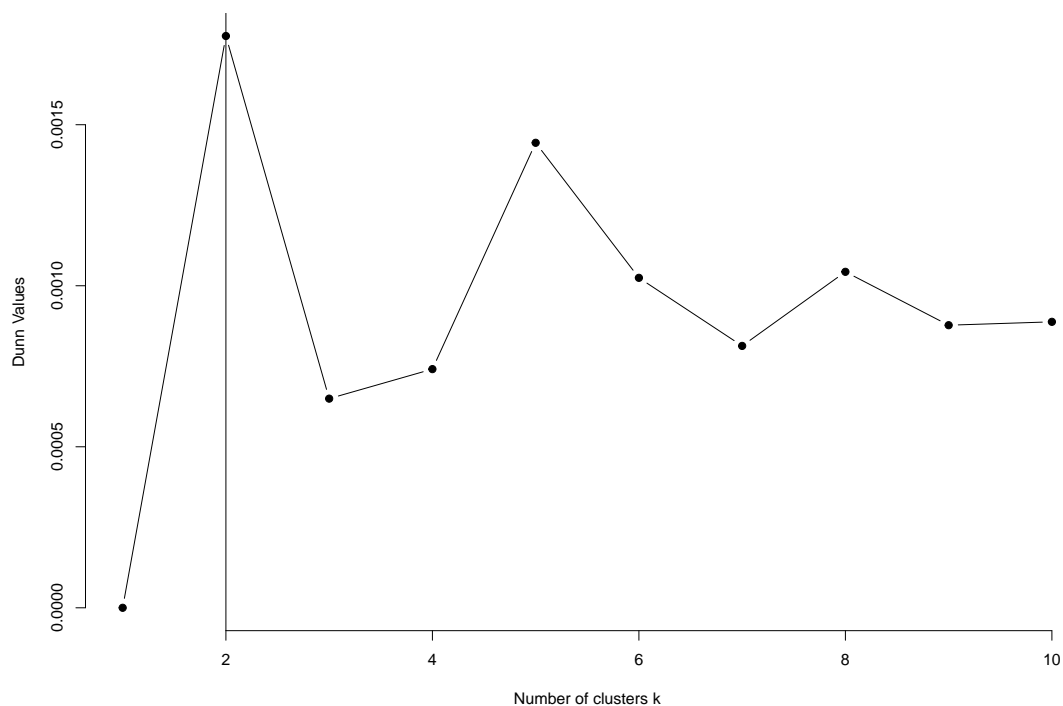


Questão 1:

Use alguma metrica interna (algum Dunn, Silhouette, Calinski-Harabaz index) - apenas uma -para escolher o k entre 2 e 10.

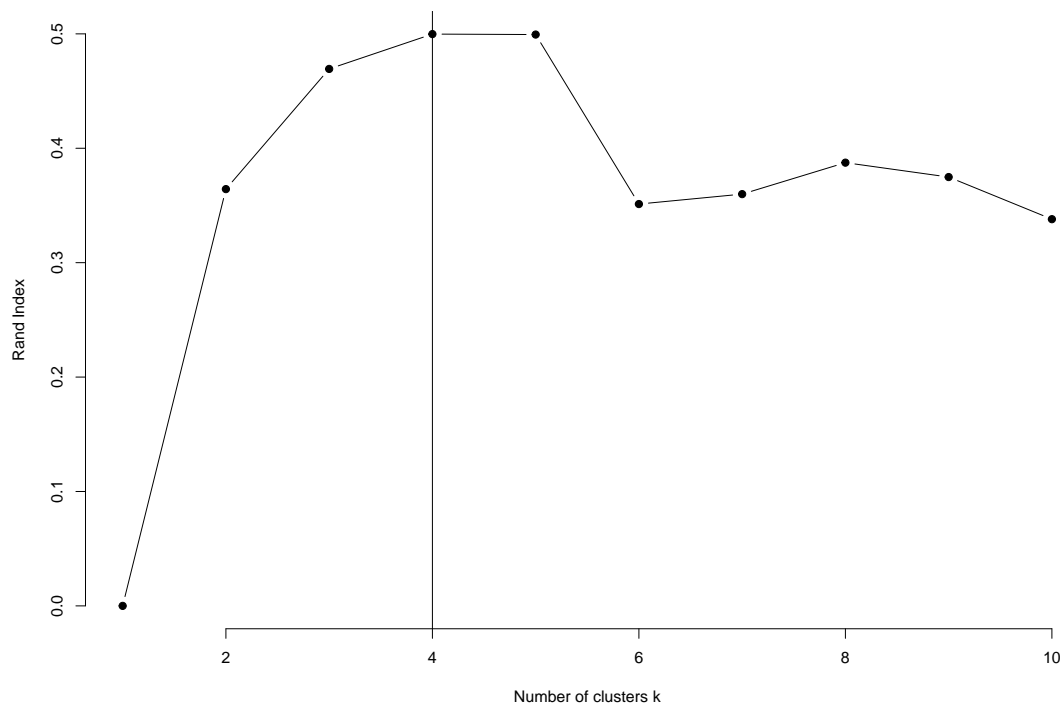
Solução: A medida interna adotada no script (ver anexo I) foi a Dunn. Ela foi extraída pela função `cluster.stats` para selecionar o valores de **k**. A comportamento dessa medida para os valores de k (de 2 até 10) é apresentado no gráfico abaixo.



Questão 2:

Use alguma medida externa (Normalized/adjusted Rand, Mutual information, variation of information) para decidir no k .

Solução: A medida externa adotada no script (ver anexo I) foi a Adjusted Rand. Ela foi extraída pela função `cluster.stats`. O valor de k selecionado através dela foi igual a **4**, conforme o gráfico abaixo.



Anexo I: Script fonte em R

```
1 # -----
2 # Description:
3 #   solutions for activity 4 (MO444)
4 #
5 # Version: 1.0
6 #
7 # Author:
8 #   Luiz Alberto , gomes.luiz@gmail.com
9 #
10 # History:
11 #   Oct 20th, 2016 started
12 #
13 # To do:
14 #   -
15 # -----
16
17 if (!require(caret))
18   install.packages(caret)
```

```
19 if (!require(cluster))
20   install.packages(cluster)
21 if (!require(fpc))
22   install.packages(fpc)
23
24 library(fpc)
25 library(caret)
26 library(cluster)
27
28 # cleans up execution environment.
29 rm(list = ls())
30
31 # sets up path to data files.
32 setwd('~\\Workspace\\doutorado\\disciplinas\\mo444b\\atividades\\4')
33
34 # -----
35 # common functions
36 # -----
37 ReadDataFile <- function(name) {
38   # Reads a data file in csv format.
39   #
40   # Args:
41   #   name: file name to be read.
42   #
43   # Returns:
44   #   the data frame with rows and columns from file.
45   #
46   result <-
47     read.csv(
48       file = paste('./data/', name, sep = ''),
49       header = TRUE,
50       sep = ',',
51     )
52   return(result)
53 }
54
55 # -----
56 # main function
57 # -----
58 main <- function() {
59   # reads raw data from files.
60   cluster.data <- ReadDataFile('cluster-data.csv')
61   cluster.data.class <- ReadDataFile('cluster-data-class.csv')
62   cluster.data.scaled <- scale(cluster.data)
63
64   set.seed(123)
65   dd <- dist(cluster.data.scaled, method = "euclidean")
66
67   # external cluster validation
68   output <- data.frame(line=character())
69   e.folds <-
70     createFolds(
71       cluster.data.class$x,
72       k = 5,
73       list = TRUE,
74       returnTrain = TRUE
75     )
76   max.rand <- 0
77   for (e in 1:5) {
78     ## internal cluster validation.
79     i.folds <-
80       createFolds(
81         cluster.data.class$x,
82         k = 3,
```

```
85     list = TRUE,
86     returnTrain = TRUE
87 )
88 max.dunn <- 0
89 for (k in 2:10) {
90     dunn <- 0
91     for (i in 1:3) {
92         kmi <- kmeans(cluster.data.scaled[i.folds[[i]], ], k, nstart = 5)
93         kmi.stats <-
94             cluster.stats(dd, kmi$cluster, silhouette = FALSE)
95         dunn <- dunn + kmi.stats$dunn
96
97         output <- rbind(output, data.frame(line=sprintf(' i = %d, k = %d, dunn = %.10f\n', i, k,
98             kmi.stats$dunn)))
99     }
100     dunn <- dunn / 3
101     if (dunn > max.dunn) {
102         max.dunn <- kmi.stats$dunn
103         max.ki <- k
104     }
105 }
106 output <- rbind(output, data.frame(line=sprintf(' max.ki = %d, max.dunn = %.10f\n', max.ki,
107     kmi.stats$dunn)))
108 ## external cluster validation
109 classes <- as.numeric(cluster.data.class$x[e.folds[[e]]])
110 kme <-
111     kmeans(cluster.data.scaled[e.folds[[e]], ], max.ki, nstart = 5)
112 kme.stats <- kme.stats <- cluster.stats(dd, classes, kme$cluster, silhouette = FALSE)
113
114 output <- rbind(output, data.frame(line=sprintf(' e = %d, max.ki = %d, max.dunn = %.10f\n',
115     e, max.ki, kme.stats$corrected.rand)))
116 if (kme.stats$corrected.rand > max.rand) {
117     max.rand <- kme.stats$corrected.rand
118     max.ke <- max.ki
119 }
120
121 ## write the final k.
122 output <- rbind(output, data.frame(line=sprintf(' max.ke = %d, max.rand = %.10f\n', max.ke,
123     max.rand)))
124 write.csv(output, file='./data/result.csv', row.names = FALSE, quote = FALSE)
125
126 }
127
128 main()
```