

MC884/MO444 - APRENDIZADO DE MÁQUINA

k-means

Erik de Godoy Perillo - RA135582

Universidade Estadual de Campinas

30 de outubro de 2016

1 Introdução

O objetivo do trabalho era explorar o algoritmo de *k-means* e suas métricas de desempenho.

1.1 Implementação

A linguagem de implementação escolhida foi o R. Todo o código utilizado no relatório encontra-se na seção 4. Ao longo do documento, linhas do código serão citadas para referência no mesmo. A função `main` (linha 18 da seção 4) executa tudo que é requisitado no enunciado, mostrando os resultados.

2 Metodologia

Os parâmetros para o *k-means* encontram-se na linha 11 da seção 4. Para métrica interna, foi selecionado o `dunn2`, que mede a razão entre a mínima dissimilaridade média entre clusters e a máxima dissimilaridade média intra-clusters. Para métrica externa, foi selecionado o *corrected rand*. A seleção dos melhores *k* para as duas métricas é feita no *loop* da linha 39 e as métricas são plotadas na linha 76.

3 Resultados

Para a métrica interna, o melhor *k* foi 2, com $dunn2 = 2.739$. Para a externa, o melhor *k* foi 4, com $rand = 0.258$. A saída da execução do código encontra-se na seção 5. Os *plots* encontram-se abaixo.

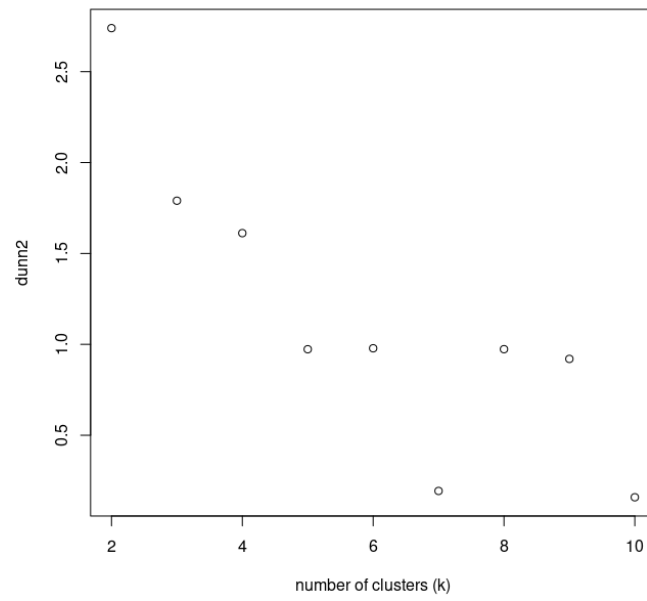


Figura 1: Métricas internas para cada k

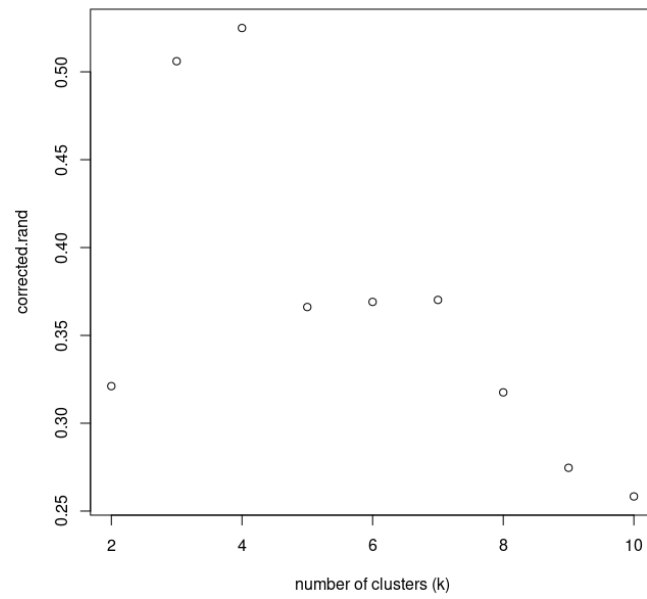


Figura 2: Métricas externas para cada k

4 Código-fonte

```
1 #packages
2 library(caret)
3 library(stats)
4 library(fpc)
5
6 #file path of data
7 data_filepath <- "cluster-data.csv"
8 labels_filepath <- "cluster-data-class.csv"
9
10 #k-means parameters
11 ks <- seq(2, 10)
12 n_start <- 5
13
14 #wrapper for sprintf
15 printf <- function(...) cat(sprintf(...))
16
17 #main method for whole challenge
18 main <- function()
19 {
20   #reading data
21   x <- read.csv(data_filepath, header=TRUE, sep=",")
22   x <- as.matrix(x)
23   y <- read.csv(labels_filepath, header=TRUE, sep=",")
24   y <- as.matrix(y[, 1])
25
26   #best index of internal metrics and external metrics
27   int_best_id <- 1
28   ext_best_id <- 1
29
30   #internal and external metrics values
31   int_metrics <- c()
32   ext_metrics <- c()
33
34   #distances between points in x
35   printf("getting distances object for points... ")
36   dst <- dist(x)
37   printf("done.\n")
38
39   for(i in seq(length(ks)))
40   {
41     printf("k = %d:\n", ks[i])
42
43     #getting k-means
44     printf("\tcomputing k-means... ")
45     means <- kmeans(x, ks[i], nstart=n_start)
46     printf("done.\n")
47
48     #getting stats
49     printf("\tcomputing clustering stats... ")
50     stats <- cluster.stats(dst, means$cluster, y)
51     printf("done.\n")
52
53     #printing metrics
54     printf("\tdunn2: %.6f | corrected.rand: %.6f\n",
55           stats$dunn2, stats$corrected.rand)
56     printf("\n")
57
58     #appending metric for later plotting
59     int_metrics <- c(int_metrics, stats$dunn2)
60     ext_metrics <- c(ext_metrics, stats$corrected.rand)
61
62     #checking for best metric
63     if(stats$dunn2 > int_metrics[int_best_id])
64       int_best_id <- i
65     if(stats$corrected.rand > ext_metrics[ext_best_id])
66       ext_best_id <- i
```

```

67     }
68
69     #printing best scores
70     printf("best internal metric score: %.6f (k=%d)\n",
71           int_metrics[int_best_id], ks[int_best_id])
72     printf("best external metric score: %.6f (k=%d)\n",
73           ext_metrics[ext_best_id], ks[ext_best_id])
74
75     #plotting metrics
76     printf("plotting metrics...\n")
77     dev.new()
78     plot(int_metrics ~ ks,
79          xlab="number of clusters (k)", ylab="dunn2")
80     dev.new()
81     plot(ext_metrics ~ ks,
82          xlab="number of clusters (k)", ylab="corrected.rand")
83
84 }

```

5 Saída do código

```
1 getting distances object for points... done.
2 k = 2:
3     computing k-means... done.
4     computing clustering stats... done.
5     dunn2: 2.739483 | corrected.rand: 0.321113
6
7 k = 3:
8     computing k-means... done.
9     computing clustering stats... done.
10    dunn2: 1.790139 | corrected.rand: 0.506041
11
12 k = 4:
13     computing k-means... done.
14     computing clustering stats... done.
15     dunn2: 1.611723 | corrected.rand: 0.524948
16
17 k = 5:
18     computing k-means... done.
19     computing clustering stats... done.
20     dunn2: 0.972936 | corrected.rand: 0.366150
21
22 k = 6:
23     computing k-means... done.
24     computing clustering stats... done.
25     dunn2: 0.972602 | corrected.rand: 0.369398
26
27 k = 7:
28     computing k-means... done.
29     computing clustering stats... done.
30     dunn2: 0.971041 | corrected.rand: 0.348167
31
32 k = 8:
33     computing k-means... done.
34     computing clustering stats... done.
35     dunn2: 0.820141 | corrected.rand: 0.290464
36
37 k = 9:
38     computing k-means... done.
39     computing clustering stats... done.
40     dunn2: 0.823787 | corrected.rand: 0.256664
41
42 k = 10:
43     computing k-means... done.
44     computing clustering stats... done.
45     dunn2: 0.155116 | corrected.rand: 0.247585
46
47 best internal metric score: 2.739483 (k=2)
48 best external metric score: 0.524948 (k=4)
49 plotting metrics...
```