

Exercício 4

Mateus Augusto Bellomo Agrello Ruivo RA:147338

28 de outubro de 2016

1 Código

```
1 from __future__ import division
2 import sklearn
3 import csv
4 import numpy as np
5 import scipy.stats as stats
6 import matplotlib.pyplot as plt
7 from numpy import genfromtxt
8 from sklearn.cluster import KMeans
9 from sklearn import metrics
10 from sklearn.metrics import pairwise_distances
11
12 X = genfromtxt('cluster-data.csv', delimiter=',', skip_header=1)
13 y = genfromtxt('cluster-data-class.csv', delimiter=',', skip_header=1)
14
15 xGraph = []
16 yGraph_adjustedRand = []
17 yGraph_adjustedMutualInfo = []
18 yGraph_vmeasure = []
19 yGraph_silhouette = []
20 yGraph_calinskiHarabaz = []
21 for k in range(2, 11):
22     xGraph.append(k)
23
24     kmeans = KMeans(n_clusters=k, n_init=5).fit(X)
25
26     yGraph_adjustedRand.append(metrics.adjusted_rand_score(y, kmeans.labels_))
27     yGraph_adjustedMutualInfo.append(metrics.adjusted_mutual_info_score(y, kmeans.labels_))
28     yGraph_vmeasure.append(metrics.v_measure_score(y, kmeans.labels_))
29
30     yGraph_silhouette.append(metrics.silhouette_score(X, kmeans.labels_, metric='euclidean'))
31     yGraph_calinskiHarabaz.append(metrics.calinski_harabaz_score(X, kmeans.labels_))
32
33
34 plt.xlabel('number of clusters (k)')
35 plt.ylabel('score value')
36 plt.plot(xGraph, yGraph_adjustedRand, 'r--', label="adjusted rand score", linewidth=2.0)
37 plt.plot(xGraph, yGraph_adjustedMutualInfo, 'b--', label="adjusted mutual info score", linewidth=2.0)
38 plt.plot(xGraph, yGraph_vmeasure, 'g--', label="v-measure score", linewidth=2.0)
39 plt.legend(bbox_to_anchor=(1, 1), bbox_transform=plt.gcf().transFigure)
40
41 plt.savefig('medidaExterna.png')
42
43
44 # normalize calinski-harabaz so we can compare it with silhouette score
45 yGraph_calinskiHarabaz = yGraph_calinskiHarabaz/np.linalg.norm(yGraph_calinskiHarabaz)
46
47 plt.gcf().clear()
48 plt.xlabel('number of clusters (k)')
49 plt.ylabel('score value')
50 plt.plot(xGraph, yGraph_silhouette, 'r--', label="silhouette coefficient", linewidth=2.0)
51 plt.plot(xGraph, yGraph_calinskiHarabaz, 'b--', label="Calinski-Harabaz index", linewidth=2.0)
52 plt.legend(bbox_to_anchor=(1, 1), bbox_transform=plt.gcf().transFigure)
53
54 plt.savefig('medidaInterna.png')
```

2 Métrica interna

O gráfico 1 compara a utilização de 2 métricas internas diferentes (*silhouette coefficient*, *Calinski-Harabaz index*) com o número de clusters variando no intervalo $[2,10]$. Pela avaliação do gráfico é possível observar que obtêm-se um melhor score para 2 clusters no método *silhouette coefficient* e para 3 clusters no método *Calinski-Harabaz index*, o que não está de acordo também com a métrica externa (descrita posteriormente).

Vale ressaltar aqui que a métrica *Calinski-Harabaz* não é normalizada. Para poder compará-la com a métrica *silhouette coefficient* no gráfico 1 realizei normalização do score.

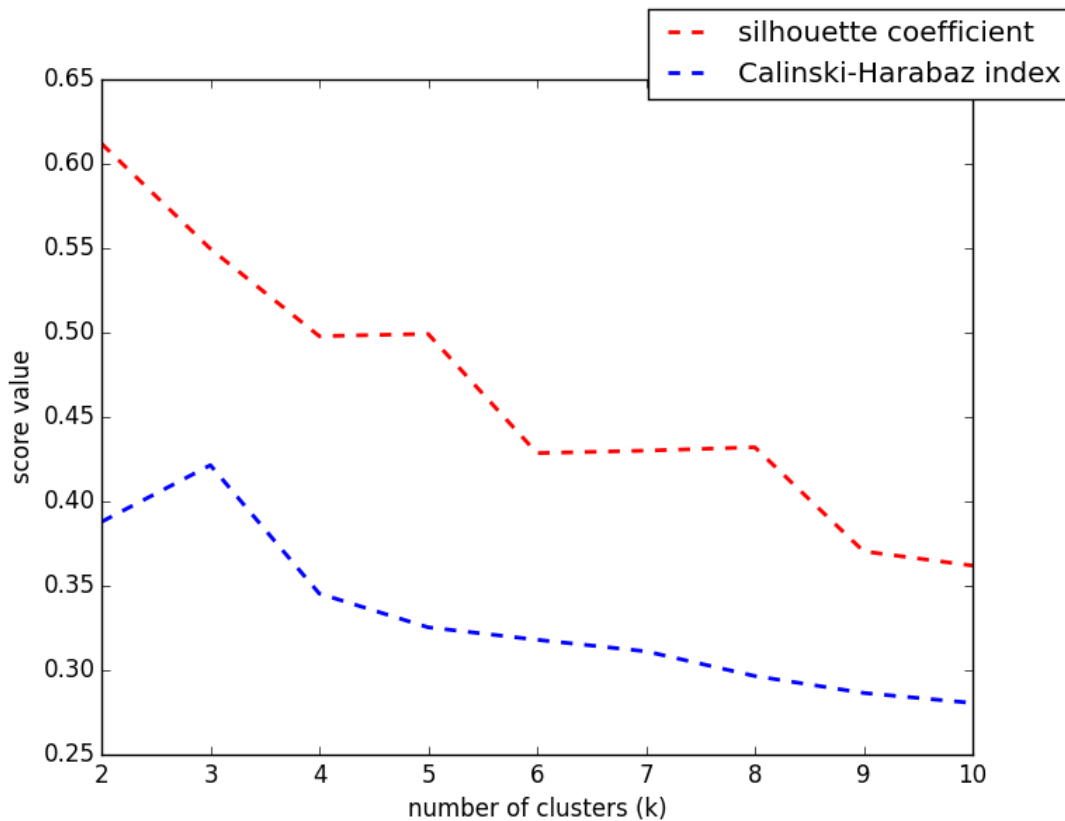


Figura 1: Métrica interna com variação de k no intervalo $[2,10]$.

3 Métrica externa

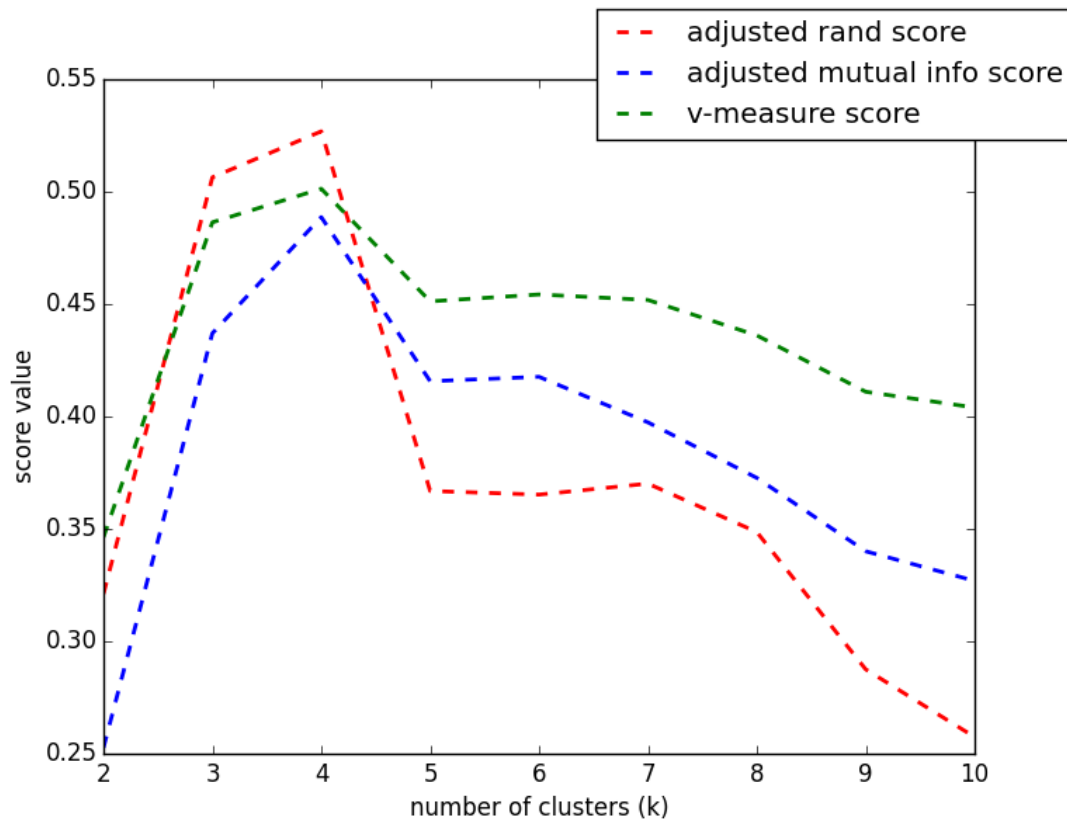


Figura 2: Métrica externa com variação de k no intervalo [2,10].

O gráfico 2 compara a utilização de 3 métricas externas diferentes (*adjusted rand score*, *adjusted mutual info score*, *v-measure score*) com o número de clusters variando no intervalo [2,10]. Pela avaliação do gráfico é possível observar que obtêm-se um melhor score para 4 clusters, e isso ocorre utilizando qualquer uma das três métricas externas.

4 Referências

- [1] <http://matplotlib.org/users/>
- [2] <http://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>
- [3] <http://scikit-learn.org/stable/modules/clustering.html#clustering-evaluation>