

Exercicio 4: KMeans

Diego Fernandez Merjildo

October 31, 2016

1 Descrição do exercicio

Use os dados do arquivo cluster-data.csv (os dados sao uma media de 30 medidas por vez da pessoa 1 do dataset Activity Recognition from Single Chest-Mounted Accelerometer Data Set

1. Rode o kmeans nos dados, com numero de restarts = 5
2. Use alguma metrica interna (algum Dunn, Silhouette, Calinski-Harabaz index) - apenas uma -para escolher o k entre 2 e 10.
3. O arquivo cluster-data-class.csv contem a classe correta de cada ponto. Use alguma medida externa (Normalized/adjusted Rand, Mutual information, variation of information) para decidir no k.
4. Plote os graficos correspondentes das 2 metricas (interna e externa) para os varios valores de k (extra).

2 Resultados

Os resultados foram obtidos rodando o script apresentado na seção 3.

Rodando o script em python obtemos os seguintes valores:

1	---	Metricas	Internas	---
2		K	Silhouette	Calinski
3		k=2	0.583	0.516
4		k=3	0.583	0.465
5		k=4	0.456	0.474
6		k=5	0.444	0.465
7		k=6	0.410	0.447
8		k=7	0.422	0.437
9		k=8	0.408	0.428
10		k=9	0.410	0.418
11		k=10	0.384	0.413

```

12 best score: 0.583031835373
13 best metric intern: Silouette
14 best k: 2
15
16 --- Metricas Externas ---
17      K    Homogen  v_meas  adj_rand  mutual
18      k=2    0.295    0.413    0.365    0.295
19      k=3    0.387    0.446    0.465    0.386
20      k=4    0.396    0.451    0.467    0.395
21      k=5    0.457    0.462    0.496    0.456
22      k=6    0.434    0.400    0.352    0.369
23      k=7    0.456    0.411    0.366    0.372
24      k=8    0.509    0.442    0.384    0.389
25      k=9    0.498    0.410    0.317    0.347
26      k=10   0.542    0.436    0.341    0.363
27 best score: 0.541596182352
28 best metric extern: Homogeneity
29 best k: 10

```

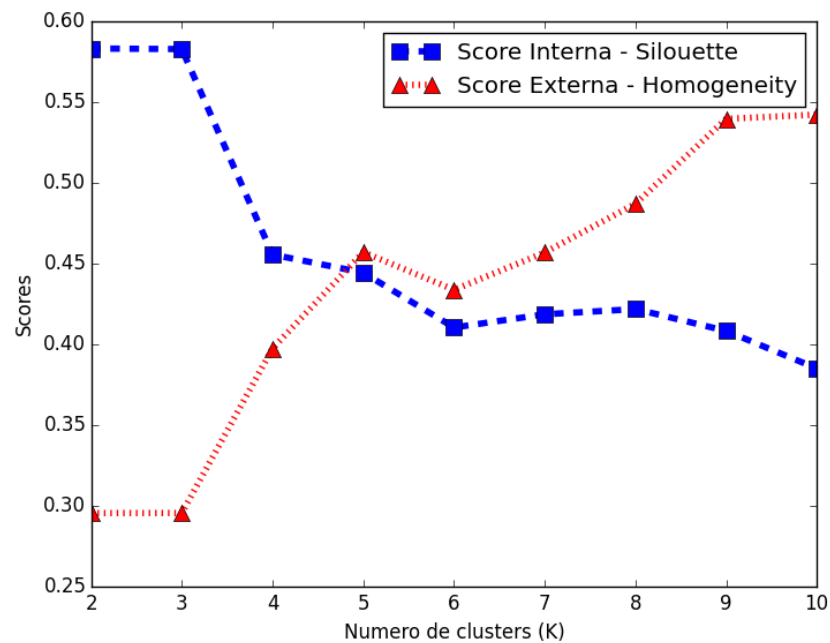


Figure 1: Figura dos scores interno e externo

3 Código fonte em python

Listing 1: Código em Python

```
1  #!/usr/bin/python
2
3  import sys,os, csv
4  import pandas
5  import numpy as np
6  import math
7  from sklearn.model_selection import StratifiedKFold
8  from sklearn.cluster import KMeans
9  from sklearn.decomposition import PCA
10 from sklearn import metrics
11 from sklearn.preprocessing import scale
12 import matplotlib.pyplot as plt
13
14 datFileName="cluster-data.csv"
15 labelsFileName="cluster-data-class.csv"
16 dirPath=os.path.dirname(os.path.realpath(__file__))
17 classList=[]
18 data=[]
19
20 def load_data(fileName):
21     raw_data = open(fileName, 'rb')
22     #rawData = pandas.read_csv(raw_data, delimiter=" ")
23     rawData = pandas.read_csv(raw_data, delimiter=",", skiprows=1)
24     return rawData.values
25
26 def getData(rawData):
27     #print "\n---- Getting data from File ----"
28     lineNum = rawData.shape[0]
29     colNum = rawData.shape[1]
30     data = np.array(rawData[0:lineNum, 0:colNum-1])
31     for i in range(lineNum):
32         classList.append(rawData[i][colNum - 1])
33     return [data, np.array(classList) ]
34
35 def get_labels(fileName):
36     labelData = load_data(dirPath + "/" + fileName)
37     labels = labelData[:,0].clip(min=0)
38     return np.array(labels)
39
40
41 def bench_k_means_inter(estimator, name, data):
42     estimator.fit(data)
43     best_score = 0.0
```

```

44     best_metric = ''
45     scores = [metrics.silhouette_score(data, estimator.labels_, ←
46         metric='euclidean', sample_size=5416),
47         metrics.calinski_harabaz_score(data, estimator.labels_ ←
48             )/10000 ]
49     print('% 9s    %.3f    %.3f'
50           % (name, scores[0], scores[1]))
51
52     if scores[0] > best_score:
53         best_score = scores[0]
54         best_metric = "Silhouette"
55
56     if scores[1] > best_score:
57         best_score = scores[1]
58         best_metric = "calinski"
59
60     return [best_score, best_metric]
61
62 def bench_k_means_ext(estimator, name, data, labels):
63     estimator.fit(data)
64     best_score = 0.0
65     best_metric = ''
66     scores = [metrics.homogeneity_score(labels, estimator.labels_),
67               metrics.v_measure_score(labels, estimator.labels_),
68               metrics.adjusted_rand_score(labels, estimator.labels_) ←
69               ,
70               metrics.adjusted_mutual_info_score(labels, estimator. ←
71               labels_) ]
72     print('% 9s    %.3f    %.3f    %.3f    %.3f'
73           % (name, scores[0], scores[1], scores[2], scores[3]))
74
75     if scores[0] > best_score:
76         best_score = scores[0]
77         best_metric = "Homogeneity"
78
79     if scores[1] > best_score:
80         best_score = scores[1]
81         best_metric = "v_measure"
82
83     if scores[2] > best_score:
84         best_score = scores[2]
85         best_metric = "adjusted_rand"
86
87     if scores[3] > best_score:
88         best_score = scores[3]
89         best_metric = "adjusted_mutual"
90
91     return [best_score, best_metric]

```

```

88
89 def main(argv=None):
90     if argv is None:
91         arv = sys.argv
92
93     data = load_data(datFileName)
94     labels = get_labels(labelsFileName)
95     data = scale(data)
96
97     best_metric_int = ''
98     best_metric_ext = ''
99     range_n_clusters = [2, 3, 4, 5, 6, 7, 8, 9, 10]
100
101     print("--- Metricas Internas ---")
102     print('%9s %s %s' % ('K', 'Silhouette', 'Calinski'))
103     best_score = 0.0
104     best_metric = ''
105     best_k = 0
106     for n_clusters in range_n_clusters:
107         new_name = "k=" + str(n_clusters)
108         [last_best_score, last_best_metric] = bench_k_means_inter(↵
            KMeans(init='k-means++', n_clusters=n_clusters, n_init↵
                =5), name=new_name, data=data)
109         if last_best_score > best_score:
110             best_score = last_best_score
111             best_metric_int = last_best_metric
112             best_k = n_clusters
113
114     print "best score: ", best_score
115     print "best metric intern: ", best_metric_int
116     print "best k:", best_k
117
118     best_score=0.0
119     print("")
120     print("--- Metricas Externas ---")
121     print('%9s %s %s %s %s' % ('K', 'Homogen', 'v_meas', '↵
        adj_rand', 'mutual'))
122     for n_clusters in range_n_clusters:
123         new_name = "k=" + str(n_clusters)
124         [last_best_score, last_best_metric] = bench_k_means_ext(↵
            KMeans(init='k-means++', n_clusters=n_clusters, n_init↵
                =5), name=new_name, data=data, labels=labels)
125         if last_best_score > best_score:
126             best_score = last_best_score
127             best_metric_ext = last_best_metric
128             best_k = n_clusters
129
130     print "best score: ", best_score

```

```

131     print "best metric extern: ", best_metric_ext
132     print "best k: ", best_k
133
134     final_score_int=[]
135     final_score_ext=[]
136
137     for n_clusters in range_n_clusters:
138         kmeans = KMeans(init='k-means++', n_clusters=n_clusters, ←
139                           n_init=5)
140         kmeans.fit(data)
141         final_score_int.append( metrics.silhouette_score(data, ←
142                                                           kmeans.labels_, metric='euclidean', sample_size=5416))
143         final_score_ext.append( metrics.homogeneity_score(labels, ←
144                                                           kmeans.labels_))
145
146     legen_int = 'Score Interna - ' + str(best_metric_int)
147     legen_ext = 'Score Externa - ' + str(best_metric_ext)
148
149     plt.plot(range_n_clusters, final_score_int, 'bs--', linewidth=4, ←
150              markersize=10, label=legen_int)
151     plt.plot(range_n_clusters, final_score_ext, 'r^:', linewidth=4, ←
152              markersize=10, label=legen_ext)
153     plt.xlabel('Numero de clusters (K)')
154     plt.ylabel('Scores')
155     plt.legend()
156     plt.show()
157
158 if __name__ == "__main__":
159     sys.exit(main())

```
